

Artificial Intelligence - Summary

École Centrale de Nantes, Davide Lanza

1. Optimal strat. in deterministic env. - 1

- Directed graph
- Strategy
- Positional strategy
- Outcome
- Objective
- Maximal sequence
- Max outcome
- Reachability of the objective
- Winning strategy
- Winning vertex
- Spinning-tree walk (BFS/DFS)
- Forgetful DFS - 2
- Completeness
- Dijkstra's algorithm
- A* algorithm

2.1 Opt. strat. in non-deterministic env. - 3

- Hypergraph
- Strategy
- Outcomes
- Solve algorithm
- Knuth extension of Dijkstra's algor. - 4

2.2 Partially observable non-det. env.

- Strategy
- Outcome
- State estimation (Filtering)
- Build a belief Hypergraph

3.1 Uncertain (probabilistic) env. - 5

- Utility of Morgenstern and V.N. (1944)
- Markov decision process
- Strategy
- Optimization in probabilistic env.
- Computing optimal strategy
- Value iteration - 6
- Policy iteration

Summary - 6.2

3.2 Partially observable deterministic env. - 7

- Partially observable MDP
- Hidden markov chains
- Probability summary reminder

- Forward algorithm
- Filtering
- Prediction
- Observation likelihood - 8
- Smoothing
- Most likely explanation (Decoding)

4.1 Supervised learning - 9

- Supervised learning
- Generalization vs. overfitting
- Loss functions
- PAC learning
- Ensemble learning

4.1 Decision trees - 10

- Build a tree
- Learning decision trees
- Finding important attributes: entropy
- Generalization and overfitting - 11
- Random forests - 12
- Decision trees: summary

4.2 Linear Models

- Linear regression
- Gradient descent
- Linear classification

4.3 Artificial neural networks - 14

- Perceptrons
- Multilayer ANN → Learning
- Gradient descent with backpropagation - 15
- Neuron saturation
- Cross-entropy & logistic function - 16
- The softmax function
- Regularisation in ANN



ART. IN. Summary

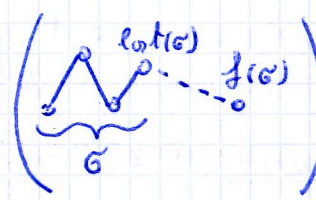
1. OPTIMAL STRATEGIES IN DETERMINISTIC ENVIRONMENTS

• DIRECTED GRAPH

A D.G. is a pair (V, E) s.t. $V =$ set of VERTICES and $E \subseteq V \times V =$ set of EDGES
(EDGE: $v \rightarrow v'$ when $(v, v') \in E$) (LABEL SETS Λ AND Σ FOR E AND V)

• STRATEGY:

A S. is a function $f: V^* \rightarrow V$ s.t. $\forall \sigma \in V^*, (\text{last}(\sigma), f(\sigma)) \in E$
(THE SET OF FINITE SEQUENCES OF ELEMENTS OF V) (LAST ELEMENT OF SEQUENCE σ)



• POSITIONAL STRATEGY

A S. f is pos. (or memoryless) if $\forall \sigma, \sigma', \text{last}(\sigma) = \text{last}(\sigma') \Rightarrow f(\sigma) = f(\sigma')$

(P.S.s are functions $V \rightarrow V$) \hookrightarrow (IS AN "ARRAY" WITH $\text{strat}(\text{state}) = \text{ACTION}$ WITH THE MDP LEXICON)

• OUTCOME

Outcome (v, f) of a strategy "f" from vertex "v" is the set of vertex sequences inductively defined by:

$$\rightarrow v \in \text{Outcome}(v, f)$$
$$\rightarrow \text{if } \sigma \in \text{Outcome}(v, f) \Rightarrow [\sigma \cdot v' \in \text{Outcome}(v, f) \Leftrightarrow v' = f(\sigma)]$$

(is the sequence of vertices produced by the agent when it apply a strategy)
(the outcome represents the successive result of each action)

$\left[\begin{array}{l} \text{ENVIRONMENT} \\ \text{FULLY OBSERVABLE} \end{array} \wedge \begin{array}{l} \text{ONLY ONE} \\ \text{AGENT} \end{array} \wedge \begin{array}{l} \text{DETERMINISTIC} \\ \text{ACTIONS} \end{array} \right] \Rightarrow \left[\begin{array}{l} \text{OUTCOME OF A STRATEGY} \\ \text{IS A SINGLE PATH} \end{array} \Rightarrow \left[\begin{array}{l} \text{THE STRATEGY CONSISTS IN } v, v' \\ \text{VERTEX OF THE PATH, GOING TO} \\ \text{THE NEXT ONE} \end{array} \right] \right]$

• OBJECTIVE

An OBS. is a set of vertex sequences

• MAXIMAL SEQUENCE

A vertex sequence is MAXIMAL within a sequence set $X \Leftrightarrow$ (is infinite) \vee \neg

\hookrightarrow (when $X = V^*$ the sequence is maximal) (is NOT A PREFIX of a sequence $\in X$)

• MAX OUTCOME

Max Outcome (v, f) is the subset of sequences of $\text{Outcome}(v, f)$ that are MAXIMAL within $\text{Outcome}(v, f)$

• REACHABILITY OF THE OBJECTIVE

Given a target vertex g , the h.o.d. $\text{Reach}(g)$ is the set of MAXIMAL seq. σ s.t. $g \in \sigma$

WINNING STRATEGY

A strategy is WINNING from some vertex "v" to some objective "X" $\Leftrightarrow \exists \text{Maximize}(v, f) \subseteq X$

WINNING VERTEX

A vertex "v" is WINNING if \exists WINNING STRATEGY from "v"

SPANNING-TREE WALK

$(s_0 \rightarrow t \text{ reachable?}) \rightarrow r = \text{false}$
 (VISITED VERTICES) $\rightarrow P = \emptyset$
 (TO-VISIT VERTICES) $\rightarrow W = s_0$

\rightarrow (if W in a queue \rightarrow BREADTH-FIRST SEARCH) ^{BFS}
 \rightarrow (if W in a stack \rightarrow DEPTH-FIRST SEARCH) ^{DFS} } PUSH-POP OPER

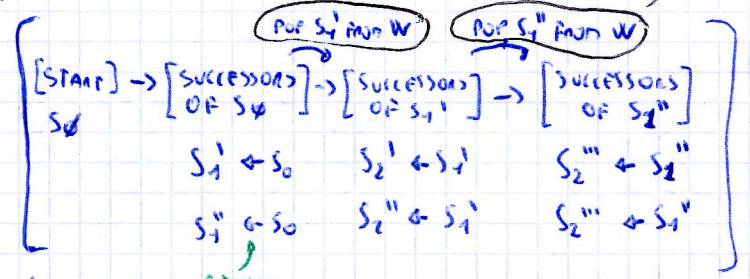
while $(W \neq \emptyset \text{ and not } r)$

POP s_i from W

if $(s_i = t) \rightarrow r = \text{true}$

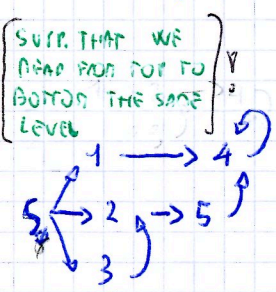
else put s_i in P AND \downarrow

$\forall s_{i+1} (\text{successor of } s_i) \notin W, P$
 PUT s_{i+1} IN W
 REGISTER s_i AS A PREDECESSOR OF s_{i+1}



REBUILD PATH TO $s_{2''''}$ $\rightarrow (s_{2''''} \leftarrow s_{1''} \leftarrow s_0)$ PATH

- EXAMPLE: $(s \rightarrow 4)$ reachable?



SO WE INSERT 1, THEN 2, THEN 3

Supp. that we read from top to bottom the same level

STEP #	0	1	2	3	4	5
W	s_0	1, 2, 3	2, 3 \rightarrow 3, 4	3, 4, 5	4, 5	5
P	\emptyset	s_0	$s_0 \rightarrow s_1$	$s_1, 2$	3, 1, 2, 3	$s_1, 2, 3$
PATH	\emptyset	s_0	$s_0 \rightarrow s_1$	$s_0 \rightarrow s_1 \rightarrow s_2$	$s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow s_3$	$s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_4$

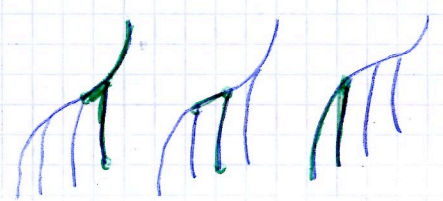
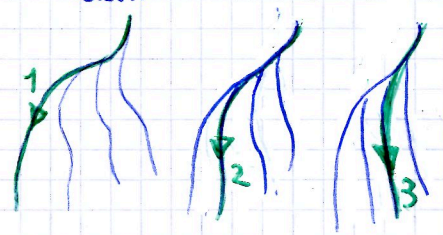
STEP #	0	1	2	3	4	5
W	s_0	1, 2, 3	1, 2	1, 5	1, 4	4
P	\emptyset	s_0	s_3	$s_4, 2, 3$	$s_0, 2, 3, 5$	$s_0, 2, 3, 5$
PATH	\emptyset	s_0	$s_0 \rightarrow s_1$	$s_0 \rightarrow s_1 \rightarrow s_2$	$s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow s_3$	$s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_4$

PL 81 "AI" BOOK

DFS \rightarrow EXPLORE A PATH UNTIL THE END, THEN THE NEAREST TO THE END etc...

F-DFS \rightarrow "FORGETFUL", REMEMBERS ONLY THE CURRENT PATH (SEE PL. AFTER)

BFS \rightarrow EXPLORE A "LEVEL" OF DEPTH, THEN FOR EACH NODE OF THE CURRENT LEVEL THE "DEEPER" LEVEL



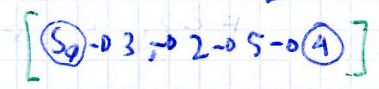
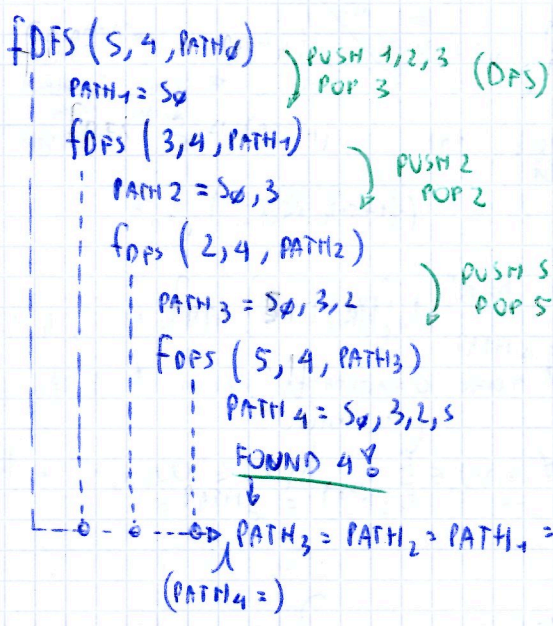
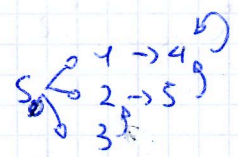
FORGETFUL DFS

~~($s_i \rightarrow t$ reachable?) $\rightarrow r = \text{false}$
 $P = \emptyset, W = \emptyset$
 while ($W \neq \emptyset$ and not r) \rightarrow NO REPLY \rightarrow
 pop s_i from W
 if ($s_i = t$) $\rightarrow r = \text{true}$
 else put s_i in P and \downarrow
 $\forall s_{i+1}$ (successor of s_i) $\notin W, P$
 put s_{i+1} in W
 register s_i on memo. of s_{i+1}~~

OUTPUT \downarrow INPUT \downarrow
 (r, PATH) function fDFS (s_i, t, PATH)
 \downarrow
 PATH' = [PATH], [s_i] \leftarrow CONCATENATION
 if ($s_i = t$) $\rightarrow r = \text{true}$
 else: $r = \text{false}$
 \forall successor s_{i+1} OF s_i AND while not r
 (r, PATH'') \leftarrow fDFS (s_{i+1}, t, PATH')
 if (r) \rightarrow PATH' = PATH'' \leftarrow RECURSIVE!

DFS/BFS \downarrow

-example:



COMPLETENESS

A search algorithm is complete if whenever \exists a path to the target vertex, it always finds one. \rightarrow (BFS complete, DFS NO). (?)

DJIKSTRA'S ALGORITHM

ALL + ∞ BUT START (R) \emptyset .
 WE WANT TO ARRIVE IN (B).
 WE CAN MOVE TO (P) AND TO (N).

85 < + ∞
 SO WE CAN MOVE TO (P)
 \hookrightarrow CHANGE THE WEIGHT TO 85
 IDEA FOR (N)

FROM (P) I CAN MOVE EVERYWHERE
 BUT:
 (R)? 85 + 85 > \emptyset NO
 (N)? 85 + 120 > 80 NO
 (B)? 85 + 124 < + ∞ YES
 \hookrightarrow I CHANGE THE WEIGHT OF (B)
 TO 85 + 124 = 209

MOVE ON P

MOVE ON N

BETTER (209)

ALL + ∞ BUT START (R) \emptyset .
 WE WANT TO ARRIVE IN (B).
 WE CAN MOVE TO (P) AND TO (N).

85 < + ∞
 SO WE CAN MOVE TO (P)
 \hookrightarrow CHANGE THE WEIGHT TO 85
 IDEA FOR (N)

FROM (P) I CAN MOVE EVERYWHERE
 BUT:
 (R)? 85 + 85 > \emptyset NO
 (N)? 85 + 120 > 80 NO
 (B)? 85 + 124 < + ∞ YES
 \hookrightarrow I CHANGE THE WEIGHT OF (B)
 TO 85 + 124 = 209

MOVE ON P

MOVE ON N

WORSE (265)

ALL + ∞ BUT START (R) \emptyset .
 WE WANT TO ARRIVE IN (B).
 WE CAN MOVE TO (P) AND TO (N).

85 < + ∞
 SO WE CAN MOVE TO (P)
 \hookrightarrow CHANGE THE WEIGHT TO 85
 IDEA FOR (N)

FROM (P) I CAN MOVE EVERYWHERE
 BUT:
 (R)? 85 + 85 > \emptyset NO
 (N)? 85 + 120 > 80 NO
 (B)? 85 + 124 < + ∞ YES
 \hookrightarrow I CHANGE THE WEIGHT OF (B)
 TO 85 + 124 = 209

MOVE ON P

MOVE ON N

BETTER (209)

ALL + ∞ BUT START (R) \emptyset .
 WE WANT TO ARRIVE IN (B).
 WE CAN MOVE TO (P) AND TO (N).

85 < + ∞
 SO WE CAN MOVE TO (P)
 \hookrightarrow CHANGE THE WEIGHT TO 85
 IDEA FOR (N)

FROM (P) I CAN MOVE EVERYWHERE
 BUT:
 (R)? 85 + 85 > \emptyset NO
 (N)? 85 + 120 > 80 NO
 (B)? 85 + 124 < + ∞ YES
 \hookrightarrow I CHANGE THE WEIGHT OF (B)
 TO 85 + 124 = 209

MOVE ON P

MOVE ON N

BETTER (209)

ALL + ∞ BUT START (R) \emptyset .
 WE WANT TO ARRIVE IN (B).
 WE CAN MOVE TO (P) AND TO (N).

85 < + ∞
 SO WE CAN MOVE TO (P)
 \hookrightarrow CHANGE THE WEIGHT TO 85
 IDEA FOR (N)

FROM (P) I CAN MOVE EVERYWHERE
 BUT:
 (R)? 85 + 85 > \emptyset NO
 (N)? 85 + 120 > 80 NO
 (B)? 85 + 124 < + ∞ YES
 \hookrightarrow I CHANGE THE WEIGHT OF (B)
 TO 85 + 124 = 209

MOVE ON P

MOVE ON N

BETTER (209)

ALL + ∞ BUT START (R) \emptyset .
 WE WANT TO ARRIVE IN (B).
 WE CAN MOVE TO (P) AND TO (N).

85 < + ∞
 SO WE CAN MOVE TO (P)
 \hookrightarrow CHANGE THE WEIGHT TO 85
 IDEA FOR (N)

FROM (P) I CAN MOVE EVERYWHERE
 BUT:
 (R)? 85 + 85 > \emptyset NO
 (N)? 85 + 120 > 80 NO
 (B)? 85 + 124 < + ∞ YES
 \hookrightarrow I CHANGE THE WEIGHT OF (B)
 TO 85 + 124 = 209

MOVE ON P

MOVE ON N

BETTER (209)

ALL + ∞ BUT START (R) \emptyset .
 WE WANT TO ARRIVE IN (B).
 WE CAN MOVE TO (P) AND TO (N).

85 < + ∞
 SO WE CAN MOVE TO (P)
 \hookrightarrow CHANGE THE WEIGHT TO 85
 IDEA FOR (N)

FROM (P) I CAN MOVE EVERYWHERE
 BUT:
 (R)? 85 + 85 > \emptyset NO
 (N)? 85 + 120 > 80 NO
 (B)? 85 + 124 < + ∞ YES
 \hookrightarrow I CHANGE THE WEIGHT OF (B)
 TO 85 + 124 = 209

MOVE ON P

MOVE ON N

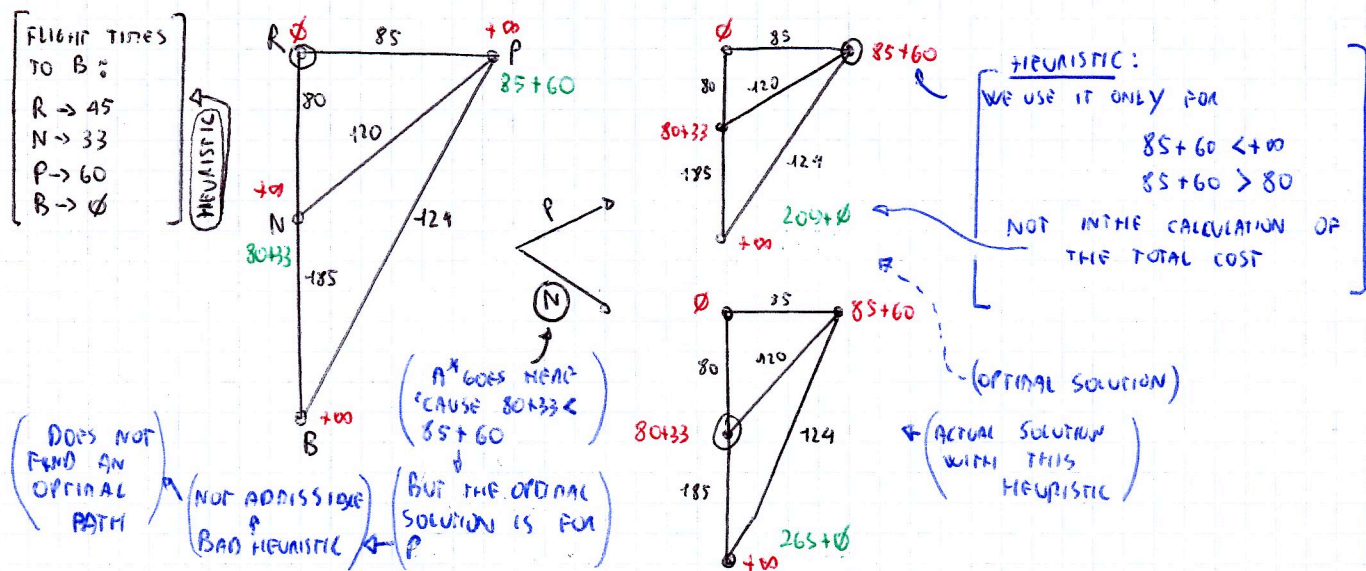
BETTER (209)

A* ALGORITHM

DIJKSTRA knows only the PATH-COST $g(m)$ TO ARRIVE TO VERTEX "m"

↳ a HEURISTIC $h(m)$ ESTIMATES the cost from vertex "m" to target "f"

DIJKSTRA minimizes $g(m)$, A* MINIMIZES $f(m) = g(m) + h(m)$



A different heuristic could be "TRAIN TIMES TO B" $\hookrightarrow [h(R)=45, h(N)=185, h(P)=124]; h(B)=∅]$ WE COULD TEST IF IT FINDS THE OPTIMAL PATH THIS TIME.

• An A* HEURISTIC is ADMISSIBLE \Leftrightarrow it never overestimates the cost of the goal

• An A* HEURISTIC is CONSISTENT $\Leftrightarrow \forall (x,y) \in E, h(x) \leq \omega(x,y) + h(y)$

• $\left. \begin{array}{l} A^* \text{ finds an OPTIMAL PATH} \leftarrow h \text{ ADMISSIBLE} \\ A^* \text{ never visits a vertex TWICE} \leftarrow h \text{ CONSISTENT} \end{array} \right\} \nabla$

(MANHATTAN DISTANCE $\Delta x + \Delta y$)
 (EUCLIDEAN DISTANCE $\sqrt{\Delta x^2 + \Delta y^2}$)

• [h CONSISTENT \Rightarrow h ADMISSIBLE]

PROOF:
 h CONSISTENT $\rightarrow \forall x, h(x) \leq c^x(x)$ with $c^x(x) = \sum_1^m \omega(x_i, x_{i+1})$

OPTIMAL PATH: (x_0, \dots, x_m)

$\left(\begin{array}{l} h(x_0) - h(x_1) \leq \omega(x_0, x_1) \\ h(x_1) - h(x_2) \leq \omega(x_1, x_2) \\ \dots \end{array} \right) \sum \rightarrow h(x_0) - \cancel{h(x_m)} \leq \sum_1^m \omega(x_i, x_{i+1})$

$(= \emptyset) \quad \checkmark \text{ QED}$

FOR THE "FLIGHT TIME" HEURISTIC:

$R \rightarrow P \rightarrow B$

$\left(\begin{array}{l} h(R) - h(P) = 45 - 60 \leq \omega(R,P) = 85 \\ h(P) - h(B) = 60 - \emptyset \leq \omega(P,B) = 124 \end{array} \right) \sum \rightarrow h(R) = 45 \leq 209$

\rightarrow [ADMISSIBLE AND CONSISTENT]

(LRTA* \rightarrow (...)) \leftarrow NOT DONE IN CLASS

2.1 OPTIMAL STRATEGIES IN NON-DETERMINISTIC ENVIRONMENTS

• HYPERGRAPH

A DIRECTED HYPERGRAPH is a pair (V, E) where $V =$ set of vertices and $E \subseteq V \times 2^V \setminus \{\emptyset\} =$ set of HYPEREDGES (ex: $(v, \{v_1, v_2\}), \dots$)
 $\underbrace{v}_S \rightarrow \underbrace{\{v_1, v_2\}}_{K'} \rightarrow$ STATE S , POSSIBLE "SUCCESSORS SET" K'

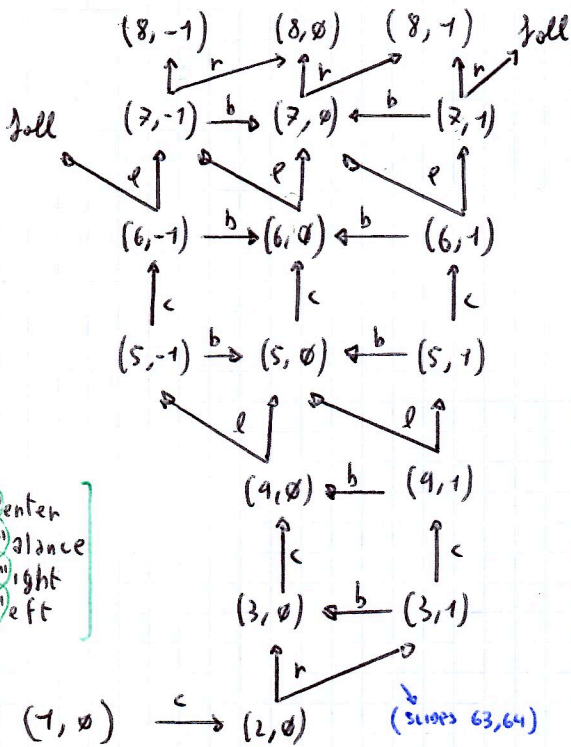
• STRATEGY

(see some on before but:) $f: V^* \rightarrow 2^V$

• OUTCOMES

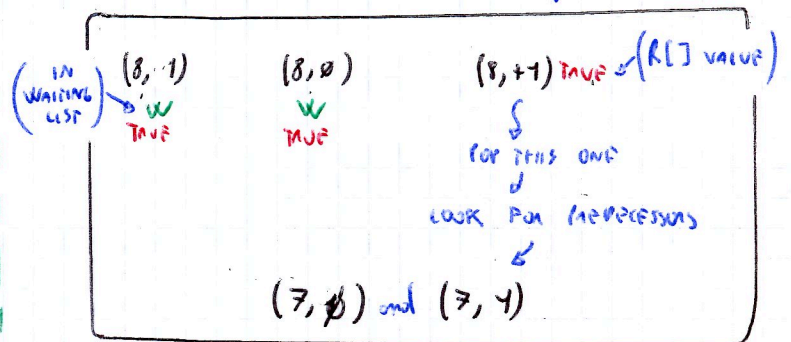
(some on before but:) a STRATEGY may have SEVERAL OUTCOMES, and the outcome is not anymore a single path but a TREE.

• SOLVE ALGORITHM

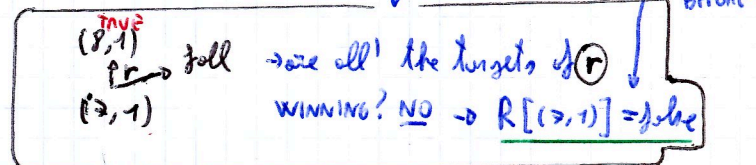


1) G is $(8, \{-1, 0, 1\}) \rightarrow$ PUT IN W

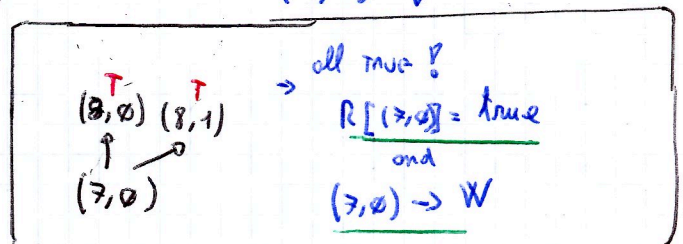
2) From my Waiting List W I'll look for a state "s"
 $\hookrightarrow (8, 1)$. My situation is \Rightarrow



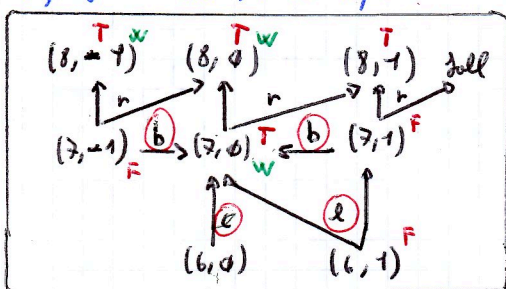
3) I start with $(7, 1) \Rightarrow$



4) Now I check $(7, 0) \Rightarrow$

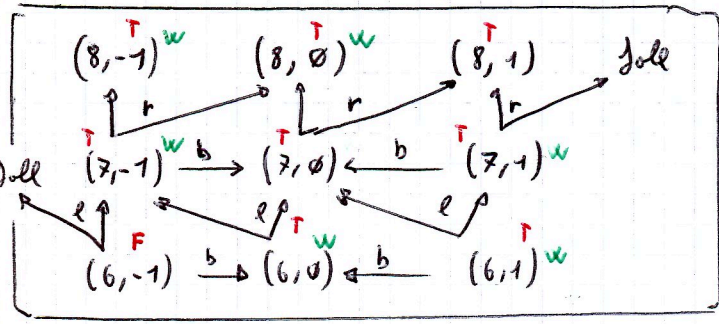


So, after these 4 steps my situation is \Rightarrow



\rightarrow I'll do the same for $(7, 0) \rightarrow$ THE WINNING ARROWS ARE HIGHLIGHTED WITH $\bigcirc \rightarrow (7, 1)$ NOW IS TRUE BECAUSE FROM $(7, 1)$ "r" IS NOT A GOOD OPTION BUT "b" YES, SO IT BECOMES TRUE AND BECOMES TRUE ALSO $(6, 1)$ WITH HIS "l" \downarrow

So, after the (3,0) ANALYSIS, we have:



→ Now, if i'll not (8,0) or (8,-1) from W, all of their predecessors are already winning → Nothing to do, continue with (7,1), (7,-1), (6,1), ...

ALGORITHM:

INPUT: s_0 // source state/root → (1,0)
 G // target state → (8, {-1,0,1})

OUTPUT: r // is G reachable from s_0 ?
 strat // the strategy

$\forall s', \text{strat}[s'] \leftarrow \perp^{\text{void symbol}}$ // No strategy (INITIALIZATION) WHEN STARTS

$r \leftarrow \text{false}$ // "is s' REACHABLE"? → ARRAY OF REACHABILITY

$\forall s' \in G, R[s'] \leftarrow \text{true}$ // all the targets are true = Reachable

$\forall s' \notin G, R[s'] \leftarrow \text{false}$ // all the other states are initially false

$W \leftarrow G$ // put the target states in the waiting list

while $W \neq \emptyset$ and not $R[s_0]$:

$s \leftarrow \text{next}(W)$ // for the next state from W (could be LIFO, FIFO, ecc...)

foreach s such that not $R[s]$ // foreach state "false" s

\hookrightarrow and $\exists (s, k) \in E$ with $s' \in K_k$ // that is a "predecessor" of s

$R[s] \leftarrow \bigvee_{(s, k') \in E} \left(\bigwedge_{s'' \in K_{k'}} R[s''] \right)$ **LIKE BEFORE**

(LOGICAL OR FOR EACH POSSIBLE SUCCESSORS SET K' OF s) (LOGICAL AND FOR EACH SUCCESSOR IN THE RESPECTIVE SUCCESSORS SET K')

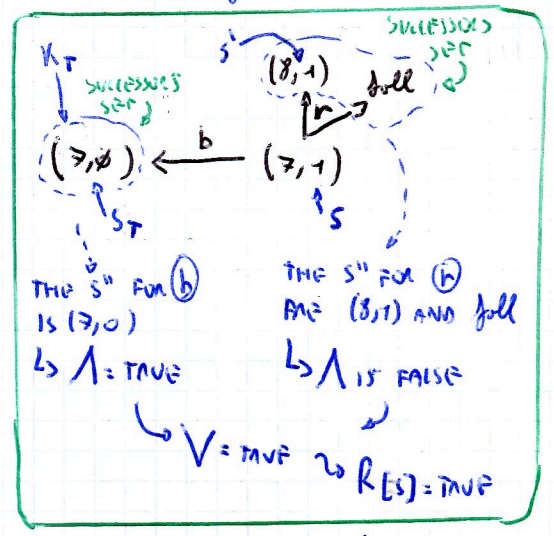
if $R[s]$:

$W \leftarrow W \cup \{s\}$ // put the true s in W

$\text{strat}[s] \leftarrow K_T$ s.t. $\bigwedge_{s_T \in K_T} R[s_T]$

PUT THE "GOOD" SUCCESSORS SET, IN THE CASE SEEN, "b" FROM (7,1), CONSTITUTED ONLY BY (7,0)

s' IS THE STATE POPPED OUT FROM W
 s IS THE NODE THAT IS TESTED
 k' IS THE INDEX FOR THE ITERATION THROUGH ALL THE SUCCESSORS SETS
 s'' IS THE INDEX FOR THE ITERATION THROUGH ALL THE STATES IN K'

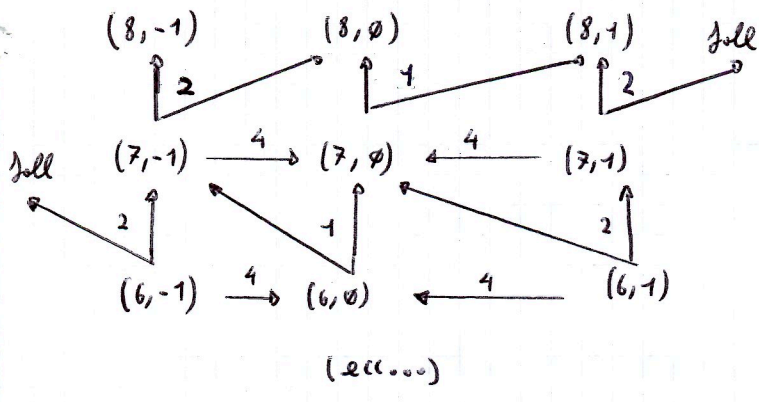


(I START WITH (8,1) → s')
 PREDECESSOR s' IS (7,1)
 I LOOP THROUGH THE K' SUCCESSORS SET AND ALL DO THE LOGICAL OPERATION.
 THEN THE ONE WHICH HAS ALL TRUE VALUES (K_T) ILL PUT IT IN THE STRATEGY

$r \leftarrow R[s_0]$

KNUTH EXTENSION OF DIJKSTRA'S ALGORITHM

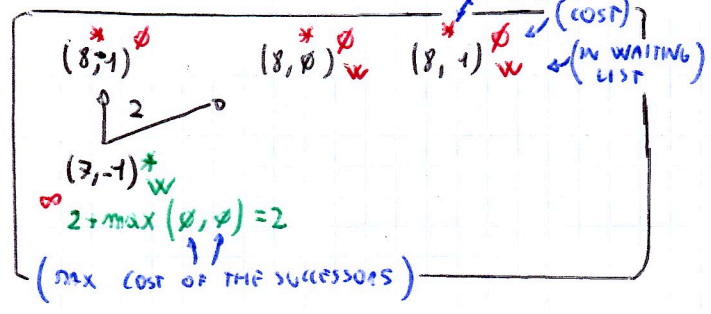
In the previous graph we'll add the costs in records of each moves: $\left[\begin{matrix} r \\ l \\ c \\ b = 4 \end{matrix} \right] \begin{matrix} 1 \text{ or } 2 \\ \\ \\ \end{matrix}$



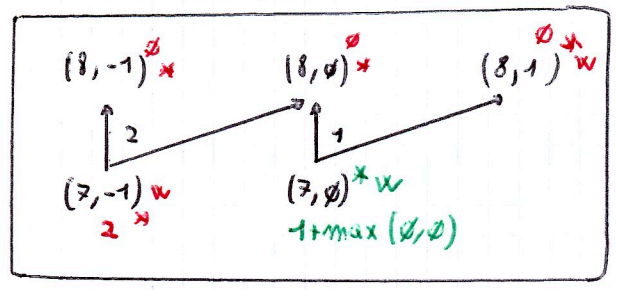
SEE SLIDE 63

The steps are:

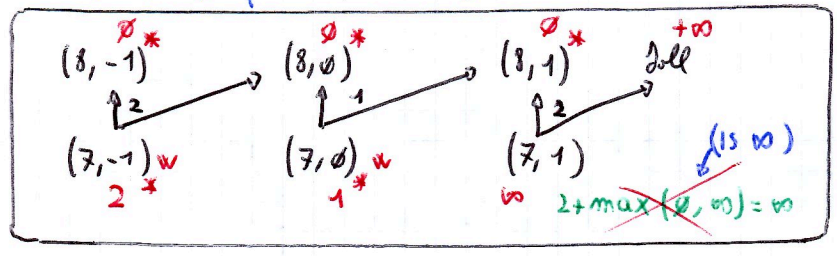
a) Will run (8,-1) from W and:



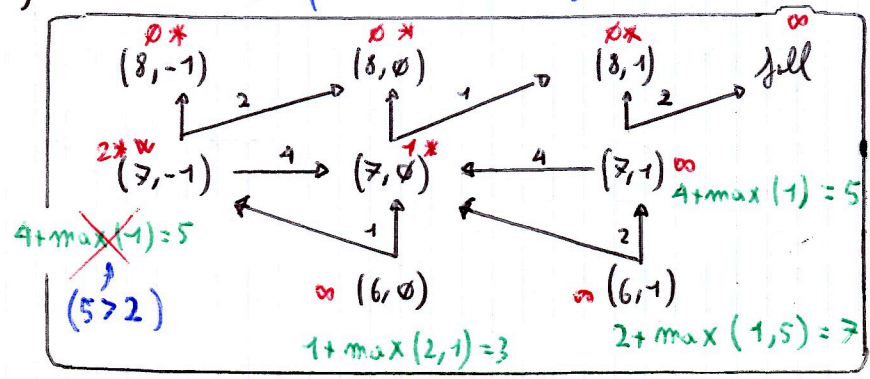
b) Will not run (7,-1) cause cost 2 > cost ∅ of the others states in W
 ↳ Will run (8,0):



c) Will run (8,1) (less cost in W):



d) Will run (7,0) (less cost in W):



ONLY FOR K' = (7,0), THE { (8,1), full } SET IS NOT CONSIDERED HERE, AS WE SEEN BEFORE

→ AND SO ON ...

2.2 PARTIALLY OBSERVABLE NON-DETERMINISTIC ENVIRONMENTS

The evolution of the system is perceived through an OBSERVATION FUNCTION mapping states to a finite set of OBSERVATIONS \mathcal{O} .

• Strategy

A strategy is a function $f: \mathcal{O}(R) \rightarrow \Sigma$ s.t. $\forall \sigma, \sigma' \in V^*$ s.t. $\mathcal{O}(\sigma) = \mathcal{O}(\sigma') \implies \exists$ one from $\text{last}(\sigma)$

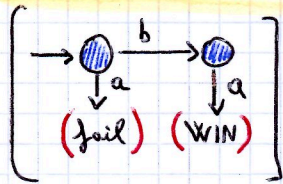
• Outcome

The Outcome (v, f) of a strategy "f" from vertex "v" is the set of runs

inductively defined by: $\rightarrow v \in \text{Outcome}(v, f)$

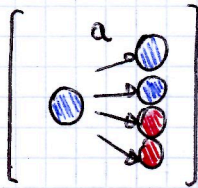
\rightarrow if $\sigma \in \text{outcome}(v, f) \implies [\sigma.v' \in \text{Outcome}(v, f) \iff v' \in f(\mathcal{O}(\sigma))]$

• State Estimation (Filtering)



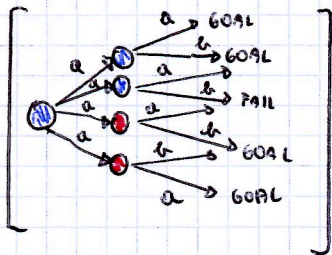
\rightarrow we observe only red/blue. Strategy: "The first time that I observe blue go (b), the second time go (a)"

\hookrightarrow NON POSITIONAL STRATEGY! (POSITIONAL: \bullet ARE ALL THE SAME, SO WE GO (a) \rightarrow FAIL)

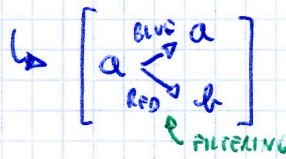


\rightarrow If I do (a) and I see blue, I FILTER the initial prediction (4 STATES) with an OBSERVATION (\rightarrow 2 STATES)

So, if we imagine a next step:



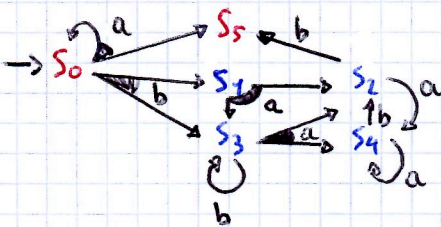
\rightarrow I'll do (a), then, if "Blue" will do (a), if "Red" will do (b)



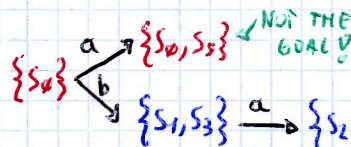
\Downarrow
We build INFORMATION SETS (BELIEF STATES)

• Build a Belief Hypergraph

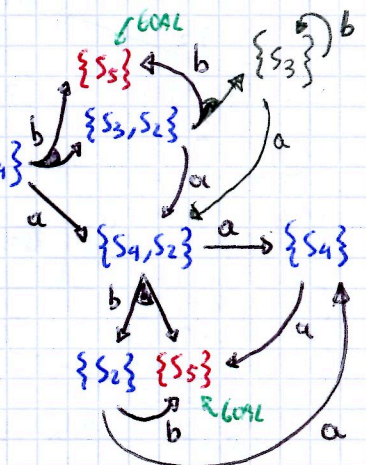
\rightarrow (SEE ALSO SLIDE 76 - COMPARISON)



BELIEF HYPERGRAPH



HERE "b" IS "NOT POSSIBLE" CAUSE FROM S4 IS POSSIBLE ONLY "a"



FINISH TO COPY PAGE

PAGE (3) EXES

3.1 UNCERTAIN (PROBABILISTIC) ENVIRONMENTS

UTILITY OF MORGENSERN AND VON NEUMANN (1921)

Preferences: (for agent d) $A \preceq B \rightarrow$ agent d weakly prefers outcome B over A

PROBABILISTIC OUTCOMES: (example: $R_1 = 10€$ if heads $\rightarrow A = 50\% \cdot R_1 + 50\% \cdot R_2 = 5$)
 $R_2 = 0€$ if tails

Utility is a function s.t.:
 1) $U(A) \leq U(B) \Leftrightarrow A \preceq B$
 2) $U(p_1 r_1 + \dots + p_n r_n) = p_1 U(r_1) + \dots + p_n U(r_n)$

example: (slide 82)

A : "You have $P = 80\%$ of winning $4K€$ "
 B : "You have $P = 100\%$ of winning $3K€$ "
 normally we prefer B ($A \preceq B$)

Is the actual value of $€$ a good utility function?

A : $0.8 \cdot 4000 + 0.2 \cdot 0 = 3200$
 B : $1 \cdot 3000 + 0 \cdot 0 = 3000$
 \rightarrow "A is better" \rightarrow $U(R_i) = R_i$ is NOT A GOOD UTILITY FUNCTION

MARKOV Decision PROCESS

A MDP is a tuple $(S, s_0, A, P, R, \gamma)$
 (FINITE SET OF STATES) (INITIAL STATE) (FINITE SET OF ACTIONS) ($P: S \times A \rightarrow \text{DIST}(S)$ IS THE TRANSITION FUNCTION) ($R: S \times A \times S \rightarrow \mathbb{R}$ IS THE REWARD FUNCTION) (DISCOUNT FACTOR) γ
 "IS AN MDP GRAPH BUT WITH PROBABILITIES AND REWARDS"

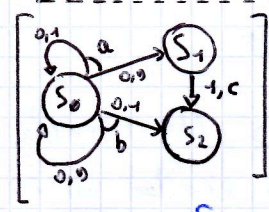
\rightarrow In the "BALANCING ROBOT" example (\rightarrow slide 85, p63)
 we just add a REWARD of 1.0 on states $\{8, \{-1, 0, 1\}\}$
 and a REWARD of -0.1 on all the other states with a 10% chance of getting embolomed \rightarrow (GRAPH AT SLIDE 86)

STRATEGY

A Strategy/Policy/Scheduler in a MDP is a function $\pi: S^* \rightarrow A$

Memoryless strategies are sufficient ONLY in a INFINITE HORIZON (slide 87)

If only two moves are accepted:



memoryless: "aa" strategy: $S_0 \xrightarrow{a} S_1 \xrightarrow{a} S_2 \rightarrow 0.9 \cdot 1$ PROBABILITY OF ARRIVING IN S_2 TOTAL: 0.9
 $S_0 \xrightarrow{a} S_0 \xrightarrow{a} S_2 \rightarrow 0$
 "bb" strategy: $S_0 \xrightarrow{b} S_2 \rightarrow 0.1$
 $S_0 \xrightarrow{b} S_0 \xrightarrow{b} S_2 \rightarrow 0.9 \cdot 0.1$ TOTAL: $0.1 + 0.09 = 0.19$

with memory: "ba" strategy: $S_0 \xrightarrow{b} S_2 \rightarrow 0.1$
 $S_0 \xrightarrow{b} S_0 \xrightarrow{a} S_2 \rightarrow 0$ TOTAL: 0.1
 "ab" strategy: $S_0 \xrightarrow{a} S_1 \xrightarrow{c} S_2 \rightarrow 0.9 \cdot 1$
 $S_0 \xrightarrow{a} S_0 \xrightarrow{c} S_2 \rightarrow 0.1 \cdot 0.1$ TOTAL: $0.9 + 0.01 = 0.91$ BEST!

FINITE HORIZONS → NONSTATIONARY OPTIMAL POLICIES
↳ WE WILL DEAL HERE WITH INFINITE HORIZON → STATIONARY OPTIMAL POLICIES

Optimization in Probabilistic Environments

No "good state" or "winning" → MAXIMISING THE REWARDS

For a single-state sequence (length = 1) → $V(s) = R_s$
(simplest case)

For sequences of states ($l > 1$) Utility is computed on the DIFFERING PART (PREFERENCES STATIONARITY)

↳ we eliminate all the common PREFIXES

$$\left[\begin{aligned} &U(s_0, s_1, \dots, s_m) < U(s_0, s_1', \dots, s_m') \\ &\Downarrow \\ &U(s_1, \dots, s_m) < U(s_1', \dots, s_m') \end{aligned} \right]$$

SEE P6.649 BOOK

So, for $\gamma \in [0, 1]$:

$$\left[U(s_0, s_1, \dots, s_m) = U(s_0) + \gamma U(s_1, \dots, s_m) \right] \rightarrow U(s_0, \dots, s_m) = \sum_{t=0}^{\infty} \gamma^t R(s_t)$$

(THE FIRST REWARD IS VERY IMPORTANT BECAUSE IT'S NOT MULTIPLIED BY γ)

With discounted rewards ($\gamma < 1$) the utility of an infinite sequence is finite:

$$\left[U(s_0, s_1, \dots, s_{\infty}) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq \sum_{t=0}^{\infty} \gamma^t R_{max} = \frac{R_{max}}{1-\gamma} \right]$$

(GEOMETRIC SERIES)

Otherwise ($\gamma = 1$) → PROPER POLICIES or AVERAGE REWARDS (see P6.650 BOOK I.A.)

Computing Optimal Strategy

Having decided that $V(s_i) = R(s_i)$ (simplest case utility = reward), we can compare policies/strategies $\{\pi\}$ by comparing the expected utilities obtained when executed.

So, given: → "s" some initial state
↳ " s_t " the state that the agent reaches at time t ($s_0 = s$) when executing π

We have:
$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \right]$$
 → Expectation is w.r.t. the probability distribution over state sequences determined by "s" and " π "

So, the optimal policy will be:
$$\pi^*(s) = \underset{\pi}{\operatorname{argmax}} [U^\pi(s)]$$

USING DISCOUNTED UTILITIES ($\gamma < 1$) WITH INFINITE HORIZONS IMPLY THAT THE OPTIMAL POLICY IS INDEPENDENT OF THE STARTING STATE (OF COURSE THE ACTION SEQUENCE WON'T BE INDEPENDENT)
↳ [SO WE CAN WRITE ONLY " π^* "] (P6.651 BOOK)

The utility function $V(s)$ allows the agent to select actions by using the principle of maximum expected utility, that is, choose the action that maximize the expected utility of the subsequent state:

$$\left[\pi^*(s) = \underset{a \in A(s)}{\operatorname{argmax}} \sum_{s'} P(s'|s, a) V(s') \right] \rightarrow \left[\begin{aligned} &2 ALGORITHMS: \\ &- VALUE ITERATION \\ &- POLICY ITERATION \end{aligned} \right]$$

VALUE ITERATION

BELLMAN EQUATION for MDP

$$V^{\pi^*}(s) = R(s) + \gamma \max_{a \in A} \sum_{s' \in S} P(s'|a,s) V^{\pi^*}(s')$$

“THE UTILITY OF A STATE IS THE IMMEDIATE REWARD FOR THAT STATE PLUS THE EXPECTED UTILITY OF THE NEXT STATE, ASSUMING THAT THE AGENT CHOOSES THE OPTIMAL ACTION”

MATIN S. 6
→ p. 652
BOOK

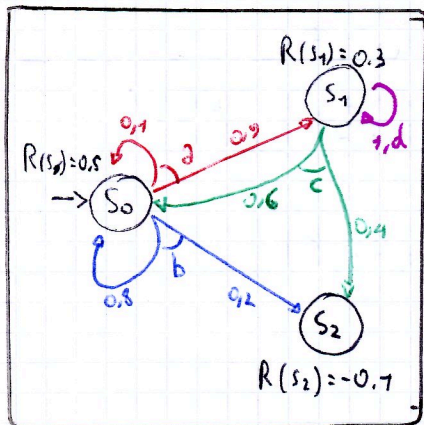
This is implied by the $\pi^*(s)$ formula seen before.

The utilities of the states ($V^{\pi^*}(s) = f[\dots]$) are solution of this θ .

These equations are non-linear, but we can solve them by iteration (up to an error ϵ):

- ALGORITHM
- $V_0(s) = \phi$
 - $V_{i+1}(s) = R(s) + \gamma \max_{a \in A} \sum_{s' \in S} P(s'|a,s) V_i(s')$ ← THIS ITERATION IS THE BELLMAN UPDATE $B(V_i)$
 - Until, $\forall s, |V_{i+1}(s) - V_i(s)| < \frac{\epsilon(1-\gamma)}{\gamma}$

↳ EXAMPLE: →



$$\begin{aligned}
 & \emptyset) \quad \underline{V_0}(s_0) = \phi \quad \underline{V_0}(s_1) = \phi \quad \underline{V_0}(s_2) = \phi \\
 & \rightarrow) \quad \underline{V_1}(s_0) = 0,5 + \gamma \cdot \max \left[P(s_0|a,s_0) V_0(s_0) + P(s_1|a,s_0) V_0(s_1), \right. \\
 & \quad \left. P(s_0|b,s_0) V_0(s_0) + P(s_2|b,s_0) V_0(s_2) \right] = 0,5 \\
 & \quad \underline{V_1}(s_1) = 0,3 \quad \underline{V_1}(s_2) = -0,1 \\
 & \quad \underline{V_1}(s_0) = 0,5 + \gamma \cdot \max \left[0,1 \cdot \phi + 0,9 \cdot \phi, 0,8 \cdot \phi + 0,2 \cdot \phi \right] = 0,5
 \end{aligned}$$

$\sum_{s' \in S} \Rightarrow$ [We get the PROBABILITY outcome for a given action "a"]
 \downarrow
 $\max_{a \in A} \Rightarrow$ [We get the action "a" that maximizes the outcome]

← HOW THE ALGORITHM WORKS
 ↳ (ITERATIONS TO MAKE) MORE ACCURACY

$$\underline{V_2}(s_0) = 0,5 + \gamma \cdot \max \left[\underbrace{0,1 \cdot 0,5 + 0,9 \cdot 0,3}_{0,32}, \underbrace{0,8 \cdot 0,5 + 0,2 \cdot (-0,1)}_{0,38} \right] = 0,88 \quad \left\{ \begin{array}{l} \text{ASSUMING} \\ \gamma = 1 \end{array} \right.$$

action "b" maximizes the outcome in S_0

$$\underline{V_2}(s_1) = 0,3 + \gamma \cdot \max \left[\underbrace{0,6 \cdot 0,5 + 0,4 \cdot (-0,1)}_{0,26}, \underbrace{1 \cdot 0,3}_{0,3} \right] = 0,6; \quad \underline{V_2}(s_2) = -0,1$$

action "d" maximizes the outcome in S_1

If we stop the iterations in $i=2 \rightarrow$ we would choose "b" in S_0 and "d" in S_1

CONVERGENCE OF VALUE ITERATION:

The Bellman update $B(U_i) = U_{i+1}$ is a CONTRACTION (by a γ factor) for the

MAX NORM $\|U\| = \max |U(s)| \rightarrow$

$$\|B(U) - B(U')\| \leq \gamma \|U - U'\|$$

The utilities are bounded by $(R_{\max}/(1-\gamma))$ so:

AFTER N ITERATIONS

$$\begin{aligned} \text{INITIAL ERROR: } \|U_0 - U^*\| &= \|U^*\| \leq R_{\max}/(1-\gamma) \\ \text{N-th ERROR: } \|U_m - U^*\| &\leq \gamma^m \|U_0 - U^*\| \leq \gamma^m (R_{\max}/(1-\gamma)) \end{aligned}$$

\hookrightarrow For if we wanna ERROR $< \epsilon \Rightarrow m = \left\lceil \frac{\log [R_{\max}/\epsilon(1-\gamma)]}{\log (1/\gamma)} \right\rceil$

(We can also show that $\|U_{i+1} - U_i\| \leq \epsilon \frac{1-\gamma}{\gamma} \Rightarrow \|U_{i+1} - U^*\| < \epsilon$)

$$\left[\|U_i - U^*\| \leq \epsilon \Rightarrow \|V^{T_i} - U^*\| \leq \epsilon \frac{1-\gamma}{\gamma} \right] \leftarrow \left(\frac{\text{CONVERGENCE OF VALUE ITERATION - POLICY LOSS}}{\text{POLICY LOSS}} \right)$$

(π_i : the strategy defined using V_i (might not be optimal))

• Policy Iteration \rightarrow (RL 656-6, 7 BOOK)

In practice, the policy becomes optimal BEFORE the value iteration converges

\downarrow ALTERNATIVE APPROACH:

- ALGORITHM
1. Fix on arbitrary policy π_i
 2. Compute the corresponding utilities U^{T_i}
 3. Compute the MEU \leftarrow strategy π_{i+1} maximizing U^{T_i} with one-step look-ahead
 4. Repeat until $\pi_{i+1} = \pi_i$

\hookrightarrow VALUE ITERATION: $B(U_i) = U_{i+1} = R(s) + \gamma \max_{a \in A(s)} \sum_{s' \in S} P(s'|a,s) U_i(s')$

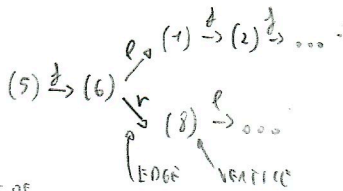
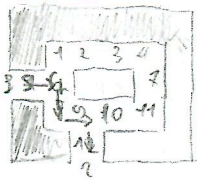
\hookrightarrow POLICY ITERATION: $B(U_i) = U_{i+1} = R(s) + \gamma \sum_{s' \in S} P(s'|a,s) U_i(s')$ \leftarrow (THIS POINTS 2. & 3.)

It simplifies the BELLMAN EQUATION! (no max \rightarrow IS A LINEAR EQUATION \rightarrow EASIER) \rightarrow (CAN BE SOLVED IN $\mathcal{O}(m^3)$ WITH LINEAR ALGEBRA)

And if we substitute step 4. with "4. Repeat for m steps" \rightarrow MODIFIED POLICY ITERATION \rightarrow (NOT EXACT BUT MORE EFFICIENT)

SUMMARY

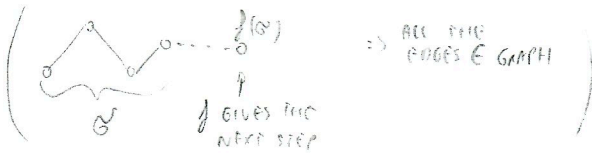
DETERMINISTIC ENVIRONMENT



Direct Graph: (V, E)
 (Edges) $\rightarrow E \subseteq V \times V$

Positional Strategy: $f: V^* \rightarrow V$
 s.t. $\forall v \in V^*, (last(\sigma), f(\sigma)) \in E$

Outcome of f = sequence of vertices reached by the agent when it applies a strategy from vertex v_0



Fully Observable: \wedge Only one default action \Rightarrow Outcome is a single path

A sequence σ is maximal in a vertex set V_0 if it is not a prefix of another sequence $\in V_0$ or it is infinite (when $V_0 = V^*$)

Of all the $\sigma \in$ outcome of f , if we took only the maximal ones σ_{max}

$$Max Outcome(v_0, f) = \{ \sigma_{max} \}$$

Objective X = is a set of vertex sequences

f is winning for $(v_0, X) \Leftrightarrow Max Outcome(v_0, f) \subseteq X$

For the game we can use as f the DFS, BFS, DFS, BFS...
 to search for a path

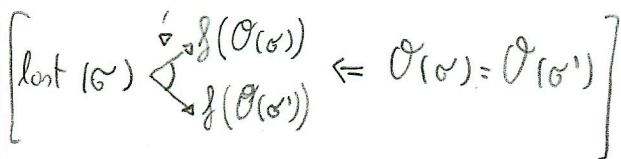
For optimal path: Dijkstra, A* \rightarrow we find

$$5 \xrightarrow{l} 6 \xrightarrow{r} 8 \xrightarrow{l} 9 \xrightarrow{r} 10$$

PARTIALLY OBSERVABLE

Observation: $O(v): V \rightarrow O_{BS}$

Strategy: $f: O(V^*) \rightarrow \Sigma$
 (set of hypotheses)



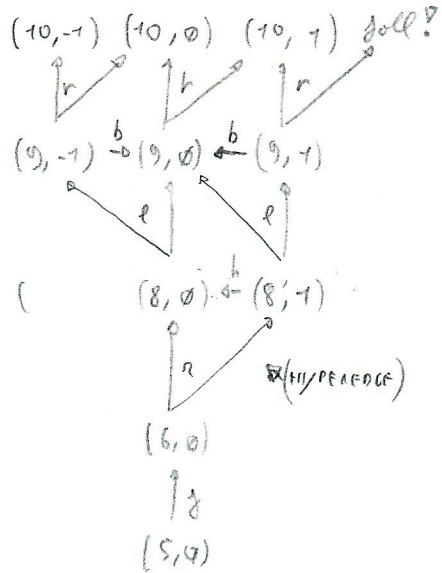
\Rightarrow [BELIEF HYPERGRAPH]

NON-DETERMINISTIC

\rightarrow we add to the game the unbalancing $(-1, \emptyset, 1)$
 left right balanced

\rightarrow we use the winning strategy for before
 $(5) \xrightarrow{r} (6) \xrightarrow{r} (8) \xrightarrow{r} (9) \xrightarrow{r} (10)$

BUT \rightarrow



Directed Hypergraph: (V, E)
 $E \subseteq V \times 2^V$ (2 edges)

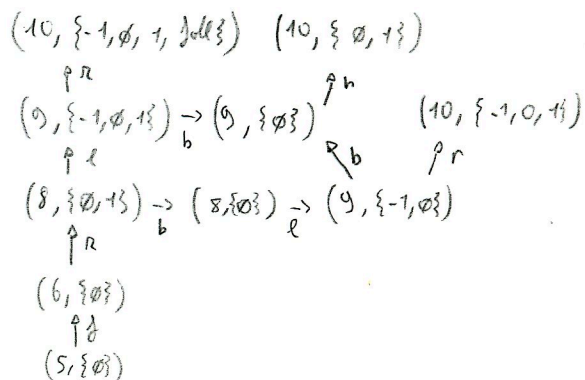
strategy: $f: V^* \rightarrow 2^V$

the outcome is a tree (the one in figure)

with "solve algorithm" I eliminate the failing possibilities from the tree \rightarrow I have again a winning strategy

if I put weights on hypergraphs

\rightarrow optimal strategy: Kuhn's extension of Dijkstra



PROBABILISTIC ENVIRONMENTS

$$\text{MDP: } (S, S_0, A, P, R^d, \gamma)$$

STATES ACTIONS REWARD
TRANSITION PROBABILITIES DISCOUNT FACTOR

↳ IS A HYPERGRAPH WITH PROBABILITIES AND REWARDS

STRATEGY \rightarrow POLICY: $\pi: S^x \rightarrow A$

NO OUTCOME \rightarrow UTILITY OF POLICY

$$V^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

↓
FIND THE OPTIMAL WITH VALUE/POLICY ITERATION

3.2 PARTIALLY OBSERVABLE PROBABILISTIC ENVIRONMENTS

• PARTIALLY OBSERVABLE MDP

OBSERVATION FUNCTION: $S \rightarrow \text{Dist}(O) \rightarrow$ ASSOCIATES to each state a prob. distribution of observations

MARKOV ASSUMPTION: $O_t \in O$ depends only on the CURRENT STATE \rightarrow we have probabilities $P(O, S)$
 \hookrightarrow difficult? \rightarrow will focus on ...

• HIDDEN MARKOV CHAINS

Po MDP where [it is possible only ONE ACTION at each time.]

HIDDEN WE WANT

- \rightarrow FILTERING: in what ^{states} state we likely to be now
- \rightarrow PREDICTION: in what states are we likely to be in the next steps
- \rightarrow OBSERV. LIKELYHOOD: how likely are the observations made?
- \rightarrow SMOOTHING: in what states were we likely to be some steps before
- \rightarrow MOST LIKELY EXPLANATION: what is the most likely sequence of states up to now?

DEFINITION OF MARKOV CHAIN

BUT BEFORE...

• PROBABILITY SUMMARY REMINDER

(ALGEBRAIC VARIABLE $X \rightarrow P(X=x)$, x "value" (or $P(x) \rightarrow P(x, y | z)$ (KNOWING THAT x, y) (MOD. OF z))
 $P(X)$ = PROBABILITY VECTOR (contains all x) $\rightarrow P(X, Y | Z)$) NOTATIONS

• PRODUCT RULE: $P(x, y) = P(x|y)P(y)$ • MARGINALIZATION (SUMMING OUT) $P(Y) = \sum_{z \in Z} P(Y, z)$

• CONDITIONING: $P(Y) = \sum_{z \in Z} P(Y|z)P(z)$

[X, Y independent] $\Rightarrow [P(X|Y) = P(X)] \wedge [P(X, Y) = P(X)P(Y)]$

[X, Y independent given Z] $\Rightarrow [P(X, Y | Z) = P(X|Z)P(Y|Z)]$

• BAYES RULE: $P(Y|X, e) = P(X|Y, e)P(Y|e) / P(X|e)$ \rightarrow (SWITCH BETWEEN CAUSALITY AND DIAGNOSTIC)
 CAUSE \rightarrow EFFECT

• FORWARD ALGORITHM

X_t is X AT TIME t (t -th step) \rightarrow OBSERVE AT EACH STEP $\Rightarrow O_{1:t} = \{O_1, \dots, O_t\}$

\hookrightarrow How to compute $P(X_t, O_{1:t})$? \rightarrow ENUMERATE ALL SEQUENCES AND SUM THE PROBABILITIES (SLOW)
FORWARD ALGORITHM

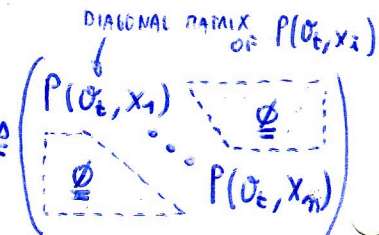
$$P(X_t, O_{1:t}) \stackrel{\text{SUMMING OUT}}{=} \sum_{x_{t-1}} P(X_t, x_{t-1}, O_{1:t}) \stackrel{\text{PRODUCT RULE}}{=} \sum_{x_{t-1}} P(O_t | X_t, x_{t-1}, O_{1:t-1}) P(X_t | x_{t-1}, O_{1:t-1}) P(x_{t-1}, O_{1:t-1})$$

(THIS IS THE OBSERVATION) MARKOV ASSUMPTION MARKOV ASSUMPTION

$\hookrightarrow P(X_t, O_{1:t}) = P(O_t | X_t) \sum_{x_{t-1}} P(X_t | x_{t-1}) P(x_{t-1}, O_{1:t-1}) \stackrel{\text{FORWARD}}{=} P(X_{t-1}, O_{1:t-1}, O_t)$
 \rightarrow ITERATIVE ALGORITHM!

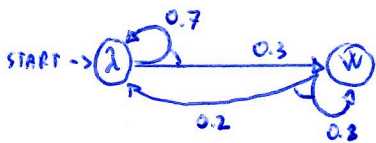
• MATRIX FORM:

$l_{1:t} \triangleq P(X_t, \sigma_{1:t})$; $T_{ij} \triangleq P(X_t = j | X_{t-1} = i) \rightarrow$ TRANSITION PROBABILITY MATRIX \underline{T}



$\hookrightarrow \underline{l}_{1:t} = \underline{\sigma}_t \cdot \underline{T}^T \cdot \underline{l}_{1:t-1} \rightarrow$ DISCRETE KALMAN FILTER

- example: (slide 113/220)



$\underline{T} = \begin{bmatrix} \lambda & w \\ 0.7 & 0.3 \\ 0.2 & 0.8 \end{bmatrix} \begin{matrix} \lambda \\ w \end{matrix}$ (3 OBSERVATIONS: S, C, P)

WHEN LISTENING (λ) STUDENTS 30% LOOK AT SCREEN (C) 10% SKIPS (S) 60% LOOK AT PROF (P)
 WHEN WEB-SURFING (W) STUDENTS 30% SKIPS (S) 60% LOOK SCREEN (C) 10% LOOK AT PROF (P)

$\underline{\sigma}_s = \begin{pmatrix} \lambda & 0.1 & \emptyset \\ w & \emptyset & 0.3 \end{pmatrix}$
 $\underline{\sigma}_c = \begin{pmatrix} \emptyset & 0.3 & \emptyset \\ \emptyset & \emptyset & 0.6 \end{pmatrix}$
 $\underline{\sigma}_p = \begin{pmatrix} \emptyset & 0.6 & \emptyset \\ \emptyset & \emptyset & 0.1 \end{pmatrix}$

we DO THE FOLLOWING OBSERVATIONS: $\sigma_{1:t} = \{S, C, P\} \leftarrow (t=3)$

$l_{1:1} = \underline{\sigma}_s T^T l_{\emptyset} \leftarrow \text{N.B. } [l_{\emptyset} = \begin{pmatrix} 1 \\ \emptyset \end{pmatrix}]^T$; $l_{1:2} = \underline{\sigma}_c T^T l_{1:1}$; $l_{1:3} = \underline{\sigma}_p T^T l_{1:2}$

\downarrow
 $\begin{pmatrix} 0.07 \\ 0.09 \end{pmatrix}$ $\begin{pmatrix} 0.021 \\ 0.058 \end{pmatrix}$ $\begin{pmatrix} 0.015 \\ 0.005 \end{pmatrix} \leftarrow \text{NORMATIVE} \rightarrow \text{OK!}$

\rightarrow TO HAVE THE PROBABILITY: NORMALIZE: $F_3 = \frac{l_{1:3}}{\text{sum}(l_{1:3})} = \begin{pmatrix} 0.75 \\ 0.25 \end{pmatrix} \rightarrow P(\lambda, S_1, C_2, P_3)$
 $\rightarrow P(W, S_1, C_2, P_3)$

FILTERING

(NOW WE CAN START)

• FILTERING

(FILTERING: in what states are we likely to be now?)

(in the last range of the example)

$P(X_t | \sigma_{1:t}) = P(X_t, \sigma_{1:t}) / P(\sigma_{1:t}) \xrightarrow{\text{MATRIX FORM}} \left[\underline{F}_{1:t} = \frac{\underline{\sigma}_t T^T \underline{l}_{1:t-1}}{\text{sum}(\underline{l}_{1:t})} = \frac{\underline{l}_{1:t}}{\text{sum}(\underline{l}_{1:t})} \right]$

• PREDICTION

(PREDICTION: in what states are we likely to be in the next items?)

$P(X_{t+m} | \sigma_{1:t}) \xrightarrow{\text{MATRIX FORM}} \left[P(X_{t+m} | \sigma_{1:t}) = (\underline{T}^T)^m \underline{F}_{1:t} \right]$

- for the example:

CONVERGES TO 40% AND 60% IF WE STOP

(We computed for S_1, C_2, P_3) \rightarrow FOR $t=13 \rightarrow \underline{F}_{1:13} = \begin{pmatrix} 0.40034 \\ 0.59966 \end{pmatrix}$, FOR $t=1003 \rightarrow \underline{F}_{1:1003} = \begin{pmatrix} 0.40000 \\ 0.60000 \end{pmatrix}$

OBSERVATION LIKELIHOOD

(how likely are the observations made?)

$$P(\sigma_{1:t}) = \sum_{x_t} P(x_t, \sigma_{1:t}) \xrightarrow{\text{marginally}} [P(\sigma_{1:t}) = \mathbf{1}^T l_{1:t}]$$

- for the example:

for $\sigma_{1:3} = \{S_1, C_2, P_3\} \rightarrow$ ASK FOR THE MOP'S CODE

SMOOTHING

(in which states we were likely to be some time before?)

We want to compute $P(X_{k:t} | \sigma_{1:t})$ (given some new evidences we refine a previous estimation)

(PROOF ON SLIDES)

$$P(X_k | \sigma_{1:t}) = \frac{1}{P(\sigma_{k+1:t})} \underline{F}_{1:k} \circ \underbrace{P(\sigma_{k+1:t} | X_k)}_{(1)}$$

$$P(\sigma_{k+1} | X_k) = \sum_{x_{k+1}} P(\sigma_{k+1} | x_{k+1}) P(\sigma_{k+2:t} | x_{k+1}) P(x_{k+1} | X_k) \quad \text{marginally}$$

$$\underline{B}_{i:t} \triangleq P(\sigma_{i:t} | X_k) \triangleq \text{BACKWARD } (\underline{B}_{i+1:t}, \sigma_i) \quad \text{with } \underline{B}_{t+1:t} = \mathbf{1}$$

↓ THE PROBABILITY OF OBSERVING AN EMPTY SEQUENCE IS ONE IN EACH STATE

$$\hookrightarrow [\underline{B}_{i:t} = \underline{T}_i \cdot \sigma_i \cdot \underline{B}_{i+1:t}] \quad (1) \text{ OK}$$

$$\text{So: } P(X_k | \sigma_{1:t}) = \frac{1}{P(\sigma_{k+1:t})} \underline{F}_{1:k} \circ \underline{B}_{k+1:t} \Rightarrow \underline{S}_{1:k:t} \triangleq \underline{F}_{1:k} \circ \underline{B}_{k+1:t}$$

$$\hookrightarrow [P(X_k | \sigma_{1:t}) = \frac{1}{\mathbf{1}^T \underline{S}_{1:k:t}} \cdot \underline{S}_{1:k:t}] \checkmark$$

- for the example:

After the filtering (compute \underline{F} vector) we want to smooth the computation

sequence: $S_1 C_2 P_3 S_4 C_5 P_6 \rightarrow$ FILTERING ON FIRST 3 $\rightarrow (k=3) \rightarrow$ done before

$$\text{So: } \underline{B}_3 = \text{ones}(2, 1) = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\underline{B}_6 = \underline{T}_6 \cdot \sigma_6 \cdot \underline{B}_3 = \begin{pmatrix} 0,45 \\ 0,2 \end{pmatrix}$$

$$\underline{B}_5 = \underline{T}_5 \cdot \sigma_5 \cdot \underline{B}_6 \quad \text{BACKWARD} = \begin{pmatrix} 0,1305 \\ 0,1230 \end{pmatrix}$$

$$\underline{B}_4 = \underline{T}_4 \cdot \sigma_4 \cdot \underline{B}_5 = \begin{pmatrix} 0,020285 \\ 0,052130 \end{pmatrix}$$

$$\underline{S}_3 = \underline{l}_3 \odot \underline{B}_4$$

"*" IN MATLAB

NORMALIZE

$$\underline{S}_3 = \underline{S}_3 / \text{sum}(\underline{S}_3) = \begin{pmatrix} 0,65263 \\ 0,34737 \end{pmatrix}$$

(SMOOTHED)	\underline{F}
↓ \underline{S}	↓
65%	75%
34%	25%

Most Likely Explanation (Decoding)

(What is the most likely sequence of states up to now?)

We want to compute: $\underset{\substack{\text{most probable sequence} \\ X_1, \dots, X_t}}{\max} P(X_{1:t} | \sigma_{1:t}) \rightarrow \left(\text{Will find the sequence of hidden states with max prob given the observations } \sigma_{1:t} \right)$

If we fix X_t , the MPS is made of:

→ a single transition from X_{t-1} to X_t

→ the most likely path to X_{t-1} given $\sigma_{1:t-1}$ (σ_t does not matter cause X_t is fixed)

→ We can compute it with VITERBI ALGORITHM

VITERBI

SCINE MS/229

4. SUPERVISED LEARNING

• SUPERVISED LEARNING

We have an UNKNOWN FUNCTION f s.t. $\forall i, y_i = f(x_i)$
 VECTOR GIVING VALUES TO SOME FINITE SET OF FEATURES OR ATTRIBUTES

So we have some given INPUT-OUTPUT PAIRS $(x_1, y_1), \dots, (x_m, y_m)$

↳ Supervised learning is finding a FUNCTION (h) (HYPOTHESIS) taken from a chosen HYPOTHESIS SPACE \mathcal{H} that APPROXIMATES f

↳ If output of f is FINITE \rightarrow CLASSIFICATION (otherwise in REGRESSION)

• GENERALISATION vs OVERFITTING

h should be consistent with given data but also generalize it \rightarrow (TESTING: K -fold cross validation)
 \hookrightarrow test data should never be used to select HYPERPARAMETERS (this would introduce bias)

[A Supervised problem is REALISABLE] $\Leftrightarrow [f(x_i) \in \mathcal{H}]$ BUT (if \mathcal{H} BIG) \rightarrow (complex learning, prob. of OVERFITTING, complex computations)

- HOW TO CHOOSE THE SIZE OF \mathcal{H} :

- 1) Let's say \mathcal{H} has functions h_i that are polynomial of degree = m
- 2) Split DATASET in TRAINING set (T_s) & VALIDATION set (V_s)
- 3) Then, For increasing values of m , train model with T_s and validate (K fold) with V_s
- 4) Choose the value of m that has the best validation, then train again with $T_s + V_s$.

• LOSS FUNCTIONS

Measures the QUALITY OF ERRORS (better than the ERROR RATE) \downarrow

$$\text{Given } y = f(x) \text{ and } \hat{y} = h(x) \Rightarrow [L(x, y, \hat{y}) \triangleq U(y, x) - U(\hat{y}, x)]$$

- EXAMPLES :

- ↳ $L_1(y, \hat{y}) = |y - \hat{y}|$ for real-valued classifiers
- ↳ $L_2(y, \hat{y}) = (y - \hat{y})^2$ " " " " "
- ↳ $L_{0/1}(y, \hat{y}) = \begin{cases} 0, & y = \hat{y} \\ 1, & y \neq \hat{y} \end{cases}$ for discrete-valued classifiers

STIPE @ UNIGR, IT
12/05 22/06 \rightarrow BUNYA (ES)

We don't minimize Expected Error Rate but the "Empirical" loss
 \hookrightarrow WE DON'T KNOW TRUE ASSOCIATED PROBABILITIES
 minimize EMPIRICAL loss

$$\text{EMPIRICAL loss: } \left[\hat{L}_E(h) = \frac{1}{N} \sum_{(x, y) \in E} L(y, h(x)) \right]$$

- REGULARIZATION OF THE LOSS FUNCTION

We add: $\left[L(\vec{x}, y, \vec{y}) \cong V(y, \vec{x}) - U(\vec{y}, \vec{x}) + \lambda \text{Reg}(h) \right]$ with $\lambda > 0$

$\text{Reg}(h) \rightarrow$ REGULARIZATION FUNCTION: is bigger when the hypothesis h is more complex

(To find the best λ): EARLY STOPPING: stop the training when performance on a VALIDATION SET first becomes worse

FEATURE SELECTION: discard statistically irrelevant attributes of the input

- PAC LEARNING ALGORITHMS

Let say that each DATA SAMPLE (\vec{x}_i, y_i) is a value of the random variable E_i

and that E_i are independent and identically distributed (I.I.D.)

$$\left\{ \begin{array}{l} \mathbb{P}(E_i | E_{i-1}, E_{i-2}, \dots) = \mathbb{P}(E_i) \\ \mathbb{P}(E_i) = \mathbb{P}(E_{i-1}) = \mathbb{P}(E_{i-2}) = \dots \end{array} \right\} \text{IID}$$

(A LEARNING ALGORITHM IS PROBABLY APPROXIMATELY CORRECT) \Leftrightarrow $\left(\exists \delta < 1, \epsilon > 0 \text{ s.t. for a SAMPLE COMPLEXITY: } \left[\mathbb{P}(\text{ERROR} \leq \epsilon) = 1 - \delta \right] \right)$

("for a big enough training set")

- EXAMPLE: BOOLEAN FUNCTION: \mathcal{Z}

$\{f: \text{BOOLEAN FUNC}\} \mathcal{Z}$ ERROR RATE: $\text{err}(h) = \mathbb{E}[L_{0,1}] \rightarrow$ (we don't know the probabilities here, so we don't know $\text{err}(h)$)
(EXPECTED VALUE)

Suppose we have a learning algorithm s.t. $\exists \epsilon > 0$ s.t. $\text{err} \leq \epsilon$

and suppose that for N -samples on hypothesis h_b is consistent but has on $\text{err}(h_b) > \epsilon$

The probability of this happening is $(1 - \epsilon)^N$ WHY??

The probability that \mathcal{H} contains such a bad hypothesis h_b is bounded by $|\mathcal{H}| (1 - \epsilon)^N$

$$\text{And } |\mathcal{H}| (1 - \epsilon)^N \leq \delta \text{ if we have: } |\mathcal{H}| (1 - \epsilon)^N \leq |\mathcal{H}| e^{-\epsilon N} \leq \delta \rightarrow \left[N \geq \frac{1}{\epsilon} \left(\log\left(\frac{1}{\delta}\right) + \log(|\mathcal{H}|) \right) \right]$$

CONDITION FOR PAC (?)

- ENSEMBLE LEARNING

Not one h but more hypothesis (K) and combine their K predictions. If the ERRORS are DISJOINT we can decrease the error a lot

- BOOSTING:

Reusing the examples in the training DATASET that are MISCLASSIFIED by the hypothesis.

They are then duplicated (or given a greater weight) to augment the original training set.

The new training set is used to obtain a new h , etc...

- BAGGING:

BOOTSTRAP AGGREGATING (\rightarrow BAGGING) consists in generating different training sets by sampling from the original training set. SAMPLING is done UNIFORMLY, and with REPLACEMENT.

Each of those new training sets (\rightarrow BOOTSTRAP SAMPLES) are used to produce hypothesis.

- RANDOM SUBSPACE:

Producing many hypothesis using only a subset of the input features for each. These subset are determined randomly for each hypothesis. The number of features can also be random, but is often the same across hypothesis.

4.1 DECISION TREES

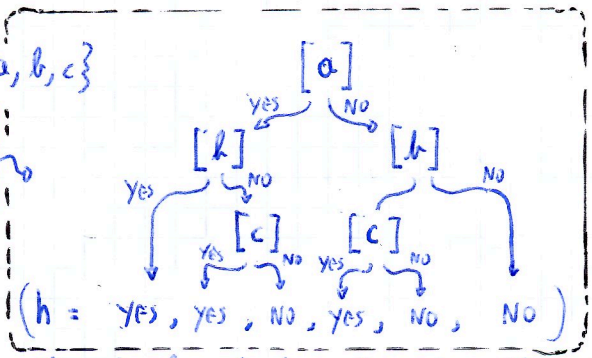
• BUILD A TREE

$(f: A \rightarrow Y \rightarrow f(\vec{x}) = y \in Y)$ $A =$ finite set of ATTRIBUTES / FEATURES

In trees the hypothesis test one ATTRIBUTE at each NODE.

- example:

here: $A = \{a, b, c\}$
 $[a, b, c, f(\vec{x})$ are
BINARY (YES/NO)



TESTS are for a finite number of values and lead to different nodes depending on the values of attributes (here: binary attributes)

LEAVES contains the final output value for the set of attributes.

(We will focus on the case where the output value is finite (CLASSIFICATION) on the possible VALUE RANGES of attributes $\in A$ are given)

Decision trees are fairly expressive (they can model exactly, for example, boolean functions: in the example, $f(\vec{x}) = (a \wedge b) \vee (a \wedge c) \vee (b \wedge c)$.) For boolean functions they can be "folded" into BDDs (BINARY DECISION DIAGRAMS, as in the example)

• LEARNING DECISION TREES

From a set of A and the corresponding values of f we wanna find the SMALLEST TREE.

We cannot use brute force (computational) \rightarrow HEURISTICS is what we need (small but not smallest trees)

THE GREEDY ALGORITHM

[EXAMPLES = VALUES OF ATTRIBUTES + CORRESPONDING VALUE OF UNKNOWN f]

- THE GREEDY ALGORITHM:

function Decision-Tree-Learning (examples, attributes, parent_examples) returns a tree

if "examples" is empty then return PLURALITY VALUE (parent_examples) // if we have no more examples, we return the plurality value (the answer that is most frequent) of the PARENT → more of the examples were matching the combination of attribute on that branch

else if all "examples" have the same classification then return their classification // try to directly ANSWER for as many examples as possible

else if "attributes" is empty then return PLURALITY VALUE (examples) // if we have no more attributes we return the plurality value of the remaining examples → there are hidden attributes to which we do not have access; or the data is inconsistent (e.g. because of noise)

else we use INFORMATION THEORY → SEE NEXT PARAGRAPH → ENTROPY

$A \leftarrow \underset{(a \in \text{attributes})}{\text{argmax}} \left[\text{IMPORTANCE}(a, \text{examples}) \right]$ // the idea is to test the most DISCRIMINATING attributes first

tree ← a new decision tree with root test A

for each value v_k of A do

exs ← { e, s.t. (e ∈ examples) ∧ (e.A = v_k) }

subtree ← Decision-Tree-Learning (exs, attributes - A, examples) for examples for which we have no direct answer we call the algorithm recursively with less example to fit, and less attribute

add a branch to tree with label (A = v_k) and subtree subtree

return tree

• FINDING IMPORTANT ATTRIBUTES: ENTROPY

Entropy of a variable V: $H(V) = - \sum_k P(v_k) \log_2(P(v_k))$ → MEASURES THE UNCERTAINTY OR V AS THE NUMBER OF BITS OF INFORMATION OBTAINED FROM A VALUE OF V

- examples:

entropy of a random variable with only one value: $H(V) = -1 \log_2(1) = 0$

entropy of a fair coin flip: $H(V) = -0.5 \log_2(0.5) - 0.5 \log_2(0.5) = -0.5(-1) + 0.5(-1) = 1$

entropy of a 4-sided die: $H(V) = -4 \cdot \frac{1}{4} \log_2\left(\frac{1}{4}\right) = 2$

The set of EXAMPLES can be seen as a random variable with: VALUES → THE OUTPUT VALUES
PROBABILITIES → " " " " FREQUENCIES

↳ We will consider only BINARY outputs (yes/no) → $P(V = \text{true}) = \frac{P}{P + M}$ ← (FREQUENCIES) } P = POSITIVE MEAN. (#YES)
M = NEG. MEAN. (#NO)

For a random boolean variable that is true with probability q, the entropy is:

$$H_2(q) = -(q \cdot \log_2(q) + (1-q) \log_2(1-q))$$

If we have an attribute $A_{\text{train}} = \{v_k \mid k \in [1, d]\}$ we'll divide the examples in d subsets E_k .

Each E_k contains p_k examples with output "yes" and m_k with output "no".

The corresponding entropy is: $H_B\left(\frac{p_k}{p_k+m_k}\right)$

The probability of having subset E_k value (Prob. that A has the k^{th} value) is $\frac{p_k+m_k}{p+m}$

So the EXPECTED ENTROPY after choosing A is:

$$R(A) = \sum_{k=1}^d \frac{p_k+m_k}{p+m} H_B\left(\frac{p_k}{p_k+m_k}\right)$$

Since we want to be as certain as possible of the value of the output after choosing A , we want to minimize this value.

INFORMATION GAIN:

minimizing $R(A) \equiv$ maximizing INFORMATION GAIN $\left[IG(A) = H_B\left(\frac{p}{p+m}\right) - R(A)\right]$

(is not suitable when A has many different values)

To correct the BIAS of the IG we need to compute:

INFORMATION GAIN RATIO: $IGR(A) = \frac{IG(A)}{IV(A)}$ \rightarrow INTRINSIC VALUE: $IV(A) = - \sum_{k=1}^d \frac{m_k+p_k}{m+p} \log_2\left(\frac{m_k+p_k}{m+p}\right)$

(is an estimation of the entropy of the random variable that gives a value for this attribute)

example:

$A_1 = \text{PROBLEMS}$	$A_2 = \text{SIZE}$	$Y = \text{DECISION}$
BIG	SMALL	Yes
NO	SMALL	No
SMALL	BIG	Yes
BIG	BIG	Yes
BIG	BIG	No
SMALL	SMALL	No
SMALL	SMALL	Yes
BIG	SMALL	Yes
NO	BIG	No
SMALL	BIG	Yes
SMALL	BIG	No
SMALL	SMALL	Yes
SMALL	BIG	No

\Rightarrow (FREQUENCIES)

PROBLEMS	SIZE
$P_{BIG} = 3$	$P_{SMALL} = 3$
$N_{BIG} = 1$	$N_{SMALL} = 4$
$P_{SMALL} = 4$	$P_{NO} = 4$
$N_{SMALL} = 3$	$N_{NO} = 2$
$P_{NO} = 0$	
$N_{NO} = 2$	

$$R(\text{problem}) = \left(\frac{4}{13} \times H_B\left(\frac{3}{4}\right)\right) + \left(\frac{7}{13} \times H_B\left(\frac{4}{7}\right)\right) + \left(\frac{2}{13} \times H_B(\emptyset)\right)$$

$$H_B\left(\frac{3}{4}\right) = -\left(\frac{3}{4} \log_2\left(\frac{3}{4}\right) + \frac{1}{4} \log_2\left(\frac{1}{4}\right)\right) \approx 0.81$$

$$R(\text{problem}) = \left(\frac{4}{13} \cdot 0.81\right) + \left(\frac{7}{13} \cdot 0.985\right) + \left(\frac{2}{13} \cdot 0\right) \approx 0.7796$$

$$R(\text{size}) = (\dots) \quad IV(\text{problem}) = -\frac{7}{13} \log_2\left(\frac{4}{13}\right) - \frac{7}{13} \log_2\left(\frac{7}{13}\right) - \frac{2}{13} \log_2\left(\frac{2}{13}\right)$$

$$IV(\text{size}) = (\dots) \approx 1.06$$

$$IGR(\text{problem}) = \frac{IG}{IV} = \frac{0.7796}{1.06} \approx 0.735$$

GENERALIZATION AND OVERFITTING

The greedy algorithm can infer PATTERNS FROM NOISE \rightarrow the algorithm will OVERFIT creating nodes for useless attributes.

(The more training data, the less ^{there} is this problem)

Decision tree PRUNING is a technique to address this problem.

- DECISION TREE χ^2 PRUNING: (PG 305 "A.I" BOOK)

We start with a full tree. We test nodes that have ONLY LEAF NODES or SUCCESSORS.

How we decide to prune or not a node? If its proportion of positives and negatives P_k, N_k calculated with $\frac{P_k}{P_k+N_k}$, is roughly similar (\approx the same) as the general $\frac{P}{m+p}$

it is an irrelevant attribute. In this case, his IG will be close to \emptyset \rightarrow IG(A) good indicator of IRRELEVANCE. How measure it? with a)

SIGNIFICANCE TEST:

1. The test begins assuming that X only underlying pattern (NULL HYPOTHESIS ASSUMPTION)
2. We calculate how much the actual data DEVIATES from the perfect absence of pattern hypothesis
3. If the degree of deviation is statistically unlikely (usually taken to mean 5% or less) then that is considered to be a good evidence for the presence of a SIGNIFICANT PATTERN, If not we accept the null hypothesis (\rightarrow WE PRUNE THE NODE)

Computation

NULL HYP. \rightarrow IG(A|A) = \emptyset for very large numbers $\nu = m+p$ ^{DOF OF SAMPLE}

How to measure the deviation Δ ? We can express the expected values with

these proportions: $\hat{P}_k = P_k \frac{P_k+N_k}{P+N}$, $\hat{N}_k = N_k \frac{N_k+P_k}{N+P}$ (ASSUMING TRIVIAL IRRELEVANCE)

$$\Delta = \sum_{k=1}^d \frac{(P_k - \hat{P}_k)^2}{P_k} + \frac{(N_k - \hat{N}_k)^2}{N_k}$$

(SEE WINCOVA TABLE)

Let's consider an discrete variable X that follows the χ^2 DISTRIBUTION with $(\nu-1)$ DOF

If NULL HYPOTHESIS true, then Δ is DISTRIBUTED as X !

example: $\nu = 4$ ^D, $\text{DOF} = 3$ $\rightarrow \chi^2 = 7.92 \rightarrow 1-\alpha = 95\%$ (FROM TABLE)

THIS MEANS THAT $P(X < 7.92) = 95\%$

So we will accept the null hyp. if $\Delta \leq \chi^2 = 7.92$ with

a confidence of $\alpha = 5\%$

$$\rightarrow \begin{cases} P(\Delta < 7.92) = 95\% = 1-\alpha \\ P(\Delta \geq 7.92) = 5\% = \alpha \end{cases}$$

So, the number of nodes pruned will be in function of the confidence α .

When we PRUNE A NODE, we replace it with a LEAF with the PLURALITY VALUE.

• RANDOM FORESTS

One on ENSEMBLE LEARNING METHOD, they consist in using BAGGING and the RANDOM SUBSPACE METHOD for the decision trees. This usually improves the accuracy of decision trees, while remaining quite fast and fairly explainable.

• Decision Trees: Summary

Decision trees are FAST, the results are EXPLAINABLE and they can be combined in random forests to gain ACCURACY.

We can use also LINEAR REGRESSION in the leaves for REGRESSION TREES ↴

4.2 LINEAR MODELS

• LINEAR REGRESSION

We consider now $\mathcal{H} = \left\{ \sum_{i=0}^m w_i x^i \mid i \in [0, m], w_i \in \mathbb{R} \forall i \right\}$ ← LINEAR FUNCTIONS SPACE

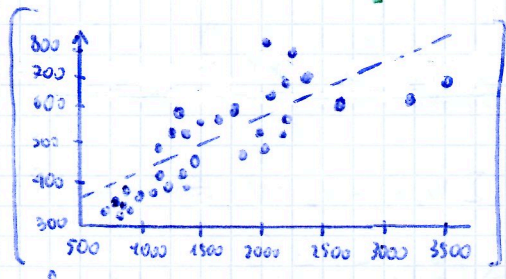
for n attributes, we have $n+1$ parameters w_i . For the REGRESSION PROBLEM we want to find the straight line $(h_{\vec{w}})$ that approximates the best set of examples.

Given told us that if the noise on data is normally distributed then this means minimizing the EUCLIDEAN DISTANCE between the actual output and the prediction.

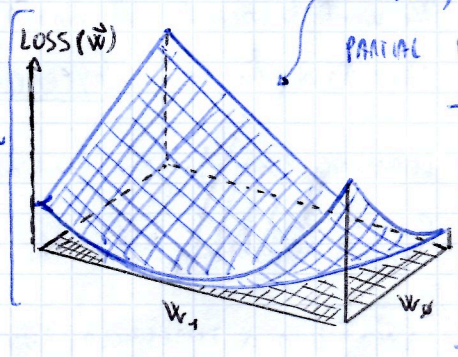
↳ minimizing L_2 over N given examples:

$$\left[\text{Loss}(h_{\vec{w}}) = \sum_{k=1}^N L_2(y_k, h_{\vec{w}}(\vec{x}_k)) \right]$$

We can find the GLOBAL MINIMUM (the function is CONVEX) by writing that all PARTIAL DERIVATIVES ARE 0



LOSS FUNCTION



REMINO: (to compute it)

$$L_2(y, \hat{y}) = (y - \hat{y})^2$$

FIGURE 15.13 (a) "IA" BOOK: data points of rice vs floor space for houses in Berkeley (CA) in 2009 → $y = (0.232)x + (246)$

(16.718 BOOK)

- OVERFITTING AND REGULARIZATION:

In dimension > 2 some attributes may exhibit spurious relations in the example set that will be captured by learning (→ OVERFITTING). To mitigate this we use REGULARIZATION on the loss function,

We consider REGULARIZATION TERMS of the form:

$$\left[\text{Loss}(h_{\vec{w}}) = \sum_{k=1}^N \left(y_k - \sum_{j=0}^m w_j x_k^j \right)^2 + \sum_{i=0}^m |w_i|^q \right]$$

REGULARIZATION L_q

Two interesting cases are $q=1$ (L_1 Regularization), that tends to produce sparse hypotheses by setting weights to 0; and $q=2$ (L_2 Regularization), more appropriate if the choice of the features is a bit arbitrary (e.g. coordinates in a rotationally invariant problem).

• Gradient Descent

With L_1 regularization $\text{Loss}(h_{\vec{w}})$ is not differentiable anymore \rightarrow instead of computing the solution in closed-form, we can use iterative methods such as

GRADIENT DESCENT:

$\vec{w} \leftarrow$ any point in the parameters space
 loop until convergence do
 for each w_i in \vec{w} do
 $w_i \leftarrow w_i - d \frac{\partial}{\partial w_i} \text{Loss}(h_{\vec{w}})$

The step size d is also called "LEARNING RATE".
 It can be fixed in time or it can decay over time as the learning process proceeds.

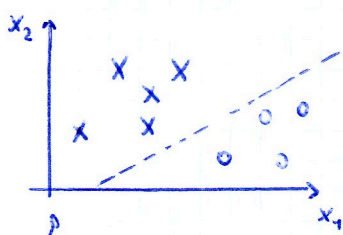
Since Loss is a big sum on a batch of examples, this is called BATCH GRADIENT DESCENT.

Convergence is guaranteed provided that d is small enough, but it will be VERY SLOW.

So, we can also randomly select one of the examples at each iteration (\rightarrow STOCHASTIC GRAD. DESC)

It is usually faster, but convergence is guaranteed only if we make d decrease at each iteration with $\left(\sum_{t=1}^{\infty} d(t) = \infty\right)$ and $\left(\sum_{t=1}^{\infty} [d(t)]^2 < \infty\right)$

• LINEAR CLASSIFICATION



(example for dimension = 2)

We want to separate all points in two categories (1 and 0). We want the separation to be on HYPERPLANE (i.e. defined by a linear equation).

LEARNING consists in finding \vec{w} s.t. $\left[\sum_i w_i x_i + w_0 \geq 0 \Leftrightarrow \text{Answer}(\vec{x}) = 1\right]$

So, our hypothesis has the form $\left[h_{\vec{w}}(\vec{x}) = T(\vec{w} \cdot \vec{x})\right]$ where

$T(\vec{w} \cdot \vec{x})$ is a THRESHOLD FUNCTION such that $T(z) = \begin{cases} 1, & z \geq 0 \\ 0, & \text{otherwise} \end{cases}$

We also assume a DUMMY COORD. x_0 that is always equal to 1.

- HARD THRESHOLD:

For $h_{\vec{w}}(\vec{x}) = T(\vec{w} \cdot \vec{x})$, $\text{Loss}(h_{\vec{w}})$ is not differentiable and the gradient is almost always 0.

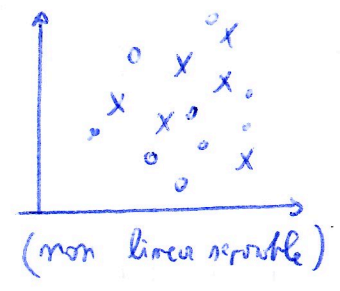
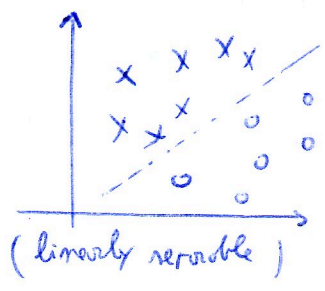
We have to find something different.

For any example (\vec{x}_k, y_k) , intuitively:

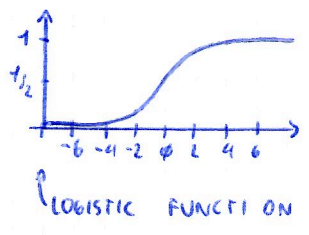
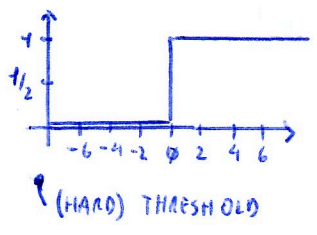
1. if $y = h_{\vec{w}}(\vec{x})$: the weights \vec{w} are fine
2. if $y=1$ and $h_{\vec{w}}(\vec{x}) = 0$: we want $h_{\vec{w}}$ to be bigger (i.e. w_i to be bigger if x_i is positive and smaller otherwise)
3. if $y=0$ and $h_{\vec{w}}(\vec{x}) = 1$: we want $h_{\vec{w}}$ to be smaller

The PERCEPTION LEARNING RULE does this: $[W_i \leftarrow W_i + d(y - h_{\vec{w}}(\vec{x}))x_i]$

[CONVERGENCE is guaranteed if] the data is LINEARLY SEPARABLE. \rightarrow Otherwise, we need to make d decrease over iterations



- LOGISTIC THRESHOLD:



Instead of a hard threshold we can use other SIMILARLY SHAPED function. The LOGISTIC FUNCTION (a sigmoid) is a good one:

$$[L(z) = \frac{1}{1 + e^{-z}}]$$

It has values between 0 and 1, so we decide by rounding to the nearest value. It is DIFFERENTIABLE. We can then compute the gradient update (using $L'(z) = L(z)(1 - L(z))$):

$$[W_i \leftarrow W_i + d(y - h_{\vec{w}}(\vec{x})) h_{\vec{w}}(\vec{x})(1 - h_{\vec{w}}(\vec{x}))x_i]$$

- EXAMPLE: LEARNING BOOLEAN FUNCTIONS

x_1	x_2	$x_1 \vee x_2$
0	0	0
0	1	1
1	0	1
1	1	1
0	0	0

Use the perception learning rule with $(d=1)$ and starting from $(w_0 = 0, w_1 = 0, w_2 = 1)$ to compute a linear classifier for $x_1 \vee x_2$ (using Hard Threshold)

I NOW: $\cancel{x_0} w_0 + x_1 w_1 + x_2 w_2 = 0 + 0 \cdot 0 + 0 \cdot 1 = 0 \xrightarrow{T(0)} = 1$ (But $x_1 \vee x_2 = 0$)

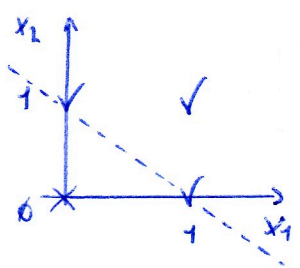
$$\text{UPDATE: } \left\{ \begin{array}{l} w_0 \leftarrow w_0 + 1(0 - 1) = -1 \\ w_1 \leftarrow w_1 + 1(0 - 1) \cdot 0 = 0 \\ w_2 \leftarrow 1 \end{array} \right.$$

II ROW: $-1 + \overset{x_1}{\emptyset} - \overset{x_2}{\emptyset} + 1 \cdot 1 = \emptyset \xrightarrow{T(\emptyset)} = 1 \quad (x_1 \vee x_2 = 1) \quad \text{OK}$

III ROW: $-1 + 1 \cdot \overset{x_1}{\emptyset} + \overset{x_2}{\emptyset} \cdot 1 = -1 \xrightarrow{T(-1)} = \emptyset \quad (x_1 \vee x_2 = 1) \rightarrow \text{UPDATE: } \{w_0 \leftarrow \emptyset, w_1 \leftarrow 1, w_2 \leftarrow 1\}$

IV ROW: $\emptyset + 1 - 1 + 1 \cdot 1 = 2 \xrightarrow{T(2)} = 1 \quad (x_1 \vee x_2 = 1) \quad \text{OK}$

V ROW: $\emptyset + \overset{x_1}{\emptyset} \cdot 1 + \overset{x_2}{\emptyset} \cdot 1 = \emptyset \xrightarrow{T(\emptyset)} = 1 \quad (x_1 \vee x_2 = \emptyset) \rightarrow \text{UPDATE: } \left. \begin{aligned} w_0 &\leftarrow w_0 + 1(\emptyset - 1) \cdot x_0 = -1 \\ w_1 &\leftarrow w_1 + 1(\emptyset - 1) \cdot x_1 = 1 \\ w_2 &\leftarrow w_2 + 1(\emptyset - 1) \cdot x_2 = 1 \end{aligned} \right\}$

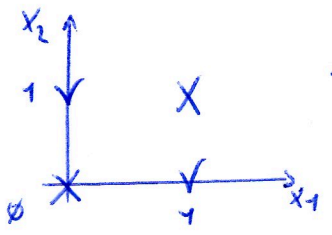


$w_0 + w_1 x_1 + w_2 x_2$ AFTER THE 5 EXAMPLES

$-1 + x_1 + x_2 \stackrel{\text{FOR PLOTTING THE LINE}}{=} \emptyset \rightarrow x_2 = -x_1 + 1$

OVER THE LINE: ALL ✓; UNDER THE LINE: ALL X \Rightarrow "LEARNING GRAPHICALLY CONFIRMED"

What for the XOR?



\rightarrow NOT LINEARLY SEPARABLE? \rightarrow no convergence with perceptron method

but with decreasing d yes?? But all...

CONCLUSION

Linear classification is often the FASTEST classification technique.

Many interesting problems are linearly separable.

Instead of computing any hyperplane, we can try to compute the one with the BIGGEST MARGIN to the nearest examples: this is the principle of Support Vector Machines (SVM), one of the most efficient classification techniques.

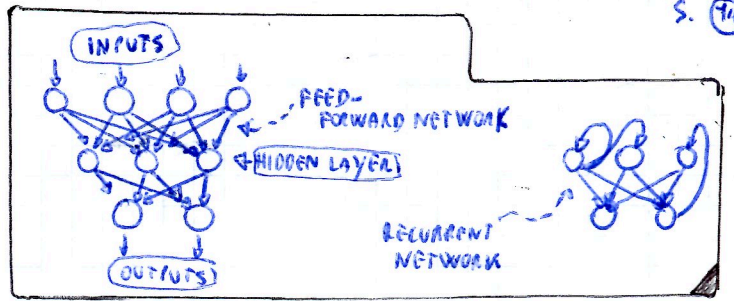
To deal with non-linearity, there are two main ideas:

- The KERNEL TRICK: consider the data in a higher dimension space. Under suitable conditions, this representation need not be explicitly computed. This is in particular applicable to SVM.

- COMBINE several linear classifiers.

4.3 ARTIFICIAL NEURAL NETWORKS

Each artificial neuron is a LINEAR CLASSIFIER.



PERCEPTRONS

A perceptron is a LINEAR CLASSIFIER

↳ ANN with a single layer (directly depends on input)

↳ each neuron has an ACTIVATION FUNCTION, it could be:

- hard threshold
 - logistic
 - $\tanh(x)$
 - identity
 - rectified linear unit: $\text{ReLU}(x) = \max(0, x)$
- } most used in SHALLOW NETWORKS
- } most used in DEEP NETWORKS

VIDEO TIP:

- Welch Labs (YT channel)
- 3Blue1Brown (YT channel)

↓

(VIDEOS ON NEURAL NETWORKS)

MULTILAYER ANN → LEARNING (MULTILAYER PERCEPTRONS)

L = # layers, σ = activation function; FULL INTERCONNECTED LAYERS (K-th input of NEURON in layer l is THE K-th NEURON OF layer $l-1$)

W_{jk}^l = weight of K-th input of J-th NEURON in L-th LAYER

z_j^l = weighted input of J-th NEURON in L-th LAYER (→ a_j^l = the output)

(Layer J dummy neuron (j=0th) s.t. $a_j^l = 1$) ⇒ So in the layer $l+1$ (VL) we can model the BIAS W_{j0}^{l+1}

So we can write:

$$z_j^l = \sum_k W_{jk}^l a_k^{l-1} = \sum_k W_{jk}^l \cdot \sigma(z_k^{l-1})$$

ACTIV. FUNC.

[let's calculate the OUTPUT LAYER ERROR]

Assuming that the loss function is ADDITIVE across the components of \vec{y} (we use L_2 here)

$$L_2(\vec{x}, \vec{y}, \hat{\vec{y}}) = \|\vec{y} - \hat{\vec{y}}\|^2 = \|\vec{y} - \vec{h}_w(\vec{x})\|^2 = \sum_k (y_k - a_k^L)^2$$

[L2 IS A PARTICULAR "LOSS" FUNCTION, WE CAN CHANGE WITH OTHERS]

(ERROR ON THE J-th NEURON) → $\left[\frac{\partial L_2}{\partial a_j^L} = -2(y_j - a_j^L) \right]$ MORE CONVENIENT → (MODIFIED ERROR) → $\left[\Delta_j^L = \frac{\partial L_2}{\partial z_j^L} = \sigma'(z_j^L) \frac{\partial L_2}{\partial a_j^L} \right]$

$\Delta_j^L \equiv \frac{\partial L_2}{\partial z_j^L}$

UPDATING WEIGHTS:

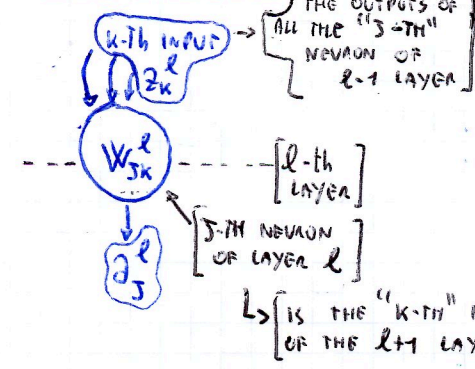
$$\frac{\partial L_2}{\partial W_{jk}^L} = \frac{\partial L_2}{\partial z_j^L} \frac{\partial z_j^L}{\partial W_{jk}^L} = \Delta_j^L \frac{\partial (\sum_i W_{ji}^L a_i^{L-1})}{\partial W_{jk}^L} = \Delta_j^L a_k^{L-1}$$

[for output layer]

$$\frac{\partial L_2}{\partial W_{jk}^l} = \Delta_j^l a_k^{l-1}$$

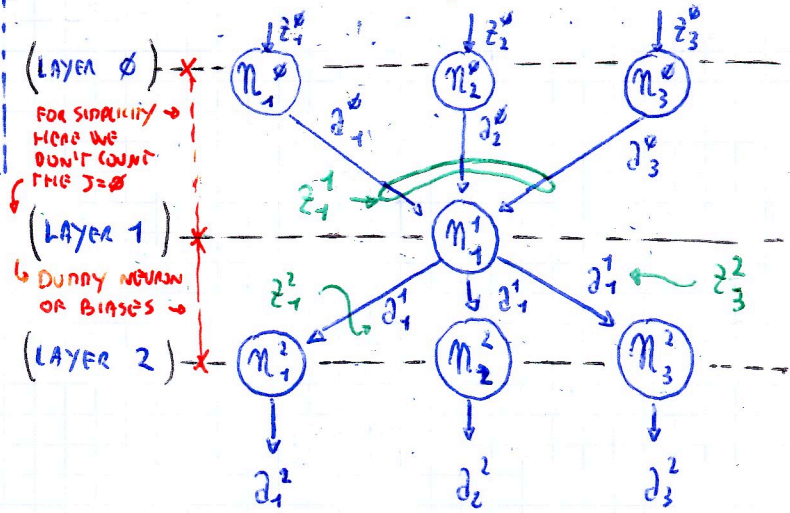
[for layer l] ← WEIGHTS ERRORS → we compute Δ_k^l with BACKPROPAGATION

BACKPROPAGATION



We got before that $\frac{\partial L_2}{\partial w_{jk}^l} = \Delta_j^l \partial_k^{l-1}$

Let's consider this example:



$n_1^1 \Rightarrow$ NEURON #1 OF LAYER 1
HAS 3 INPUTS: z_1^1, z_2^1, z_3^1
 $z_1^1 = w_{11}^1 d_1^0 + w_{12}^1 d_2^0 + w_{13}^1 d_3^0$

We compute Δ_j^l of each layer starting from the output layer error

BACKPROPAGATION OF Δ_j^l

- 1) COMPUTE OUTPUT LAYER ERROR: $\Delta_1^2 = -2(y_1 - d_1^2); \Delta_2^2 = \dots; \Delta_3^2 = \dots$
- 2) BACKPROPAGATE TO LAYER 1: $\Delta_1^1 = \sum_m g'(z_1^1) \cdot w_{m1}^2 \Delta_m^2$
 (IS NOT d_1^1 !)
 LOOP THROUGH THE LAYER 2
 # NEURON OF LAYER 2
 # INPUT OF LAYER 2 (# NEURON OF LAYER 1)
- 3) DO THE SAME FOR LAYER 0: $\Delta_1^0 = \dots; \Delta_2^0 = \dots; \Delta_3^0 = \dots$

$\Delta_{j_2}^l = \frac{\partial L_2}{\partial z_{j_2}^l} = \sum_{j_{l+1}} \frac{\partial L_2}{\partial z_{j_{l+1}}^{l+1}} \cdot \frac{\partial z_{j_{l+1}}^{l+1}}{\partial z_{j_2}^l} = \sum_{j_{l+1}} g'(z_{j_2}^l) w_{j_{l+1} j_2}^{l+1} \Delta_{j_{l+1}}^{l+1}$

(# NEURON OF LAYER l) (LOOP THROUGH THE NEURONS OF $l+1$ LAYER) (J_{l+1}TH NEURON OF LAYER $l+1$) (J_{l+1}TH INPUT FROM J_{l+1}TH NEURON OF LAYER l)

MATRIX FORM:

$\underline{w}^l = \{w_{jk}^l, v_{j,k}\}; \underline{z}^l = \{z_j^l, v_j\}; \underline{d}^l = \{d_j^l, v_j\}; \underline{\Delta}^l = \{\Delta_j^l, v_j\}$ (MATRICES FOR LAYER l)
 $\underline{\nabla}_2 L_2 = \left\{ \frac{\partial}{\partial z_j^l} L, v_j \right\}$

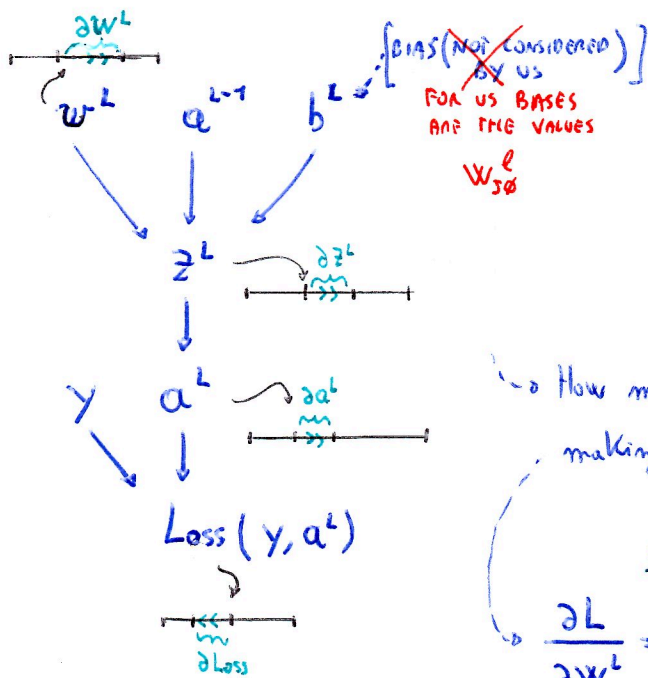
(OUTPUT LAYER ERROR) $\rightarrow \underline{\Delta}^l = \underline{\nabla}_2 L_2 \odot g'(z^l)$ (INPUT FOR THE FEED FORWARD) $\rightarrow \underline{w}^l \odot \underline{d}^{l-1} = \underline{z}^l$

(BACKPROPAGATION) $\rightarrow \underline{\Delta}^{l-1} = (\underline{w}^l)^T \underline{\Delta}^l \odot g'(z^{l-1})$

(WEIGHT UPDATE) $\rightarrow \underline{w}^l \leftarrow \underline{w}^l - \eta \underline{\Delta}^l (\underline{d}^{l-1})^T$

GRADIENT STEP (LEARNING RATE) η

USE BACKPROPAGATION TO DO GRADIENT DESCENT AS IN THE LINEAR CASE
 (CONTINUE AFTER THE "VIDEO NOTES")



If Loss = L2 \rightarrow loss = $(y - a^L)^2$

$$\begin{cases} z^L = w^L a^{L-1} + b^L & (1) \\ a^L = g(z^L) & (2) \end{cases}$$

How much a variation on w^L (∂w^L) influences Loss making it vary of ∂ Loss?

CHAIN RULE

$$\frac{\partial L}{\partial w^L} = \frac{\partial z^L}{\partial w^L} \cdot \frac{\partial a^L}{\partial z^L} \cdot \frac{\partial L}{\partial a^L}$$

(3) (2) (1)

- \hookrightarrow (1) $\frac{\partial L}{\partial a^L} = 2(y - a^L)$
 - \hookrightarrow (2) $\frac{\partial a^L}{\partial z^L} = g'(z^L)$
 - \hookrightarrow (3) $\frac{\partial z^L}{\partial w^L} = a^{L-1}$
- (1) * (2) = Δ^L FOR US

NOTE: How much a variation of BIAS b^L (∂b^L) influences ∂ Loss? (WE DON'T SEE BIAS WE USE w_{j0}^L)

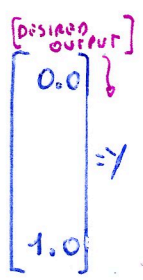
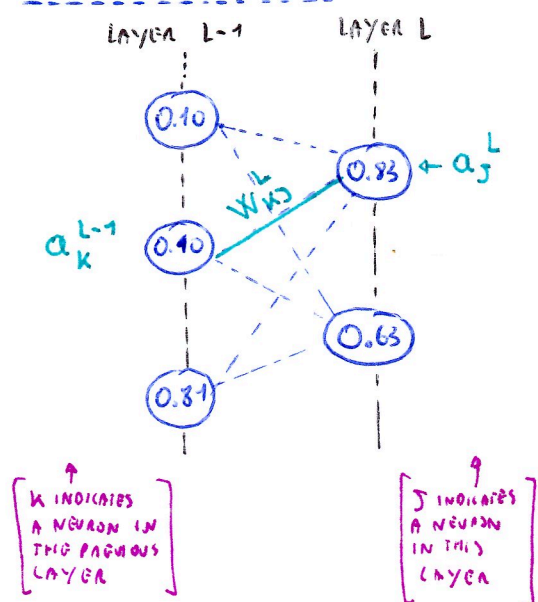
(1) & (2) AS FOR w^L

(3) $\frac{\partial z^L}{\partial b^L} = 1$ \leftarrow AS WE CAN SEE FROM EQUATION (1)

NOTE: How much a variation of a^{L-1} (∂a^{L-1}) influences ∂ Loss? \rightarrow (1) & (2) SAME FOR w^L

(3) $\frac{\partial z^L}{\partial a^{L-1}} = w^L$

- GENERALIZE THE NOTATION



\rightarrow Loss = $\sum_{j=0}^{m_L-1} (y_j - a_j^L)^2$

\rightarrow $z_j^L = \sum_{k=0}^{m_{L-1}} (w_{jk}^L a_k^{L-1}) + b_j^L$

\rightarrow $a_j^L = g(z_j^L)$

CHAIN RULE: $\frac{\partial L}{\partial w_{jk}^L} = \frac{\partial z_j^L}{\partial w_{jk}^L} \cdot \frac{\partial a_j^L}{\partial z_j^L} \cdot \frac{\partial Loss}{\partial a_j^L}$

BUT \rightarrow

$$\frac{\partial Loss}{\partial a_k^{L-1}} = \sum_{j=0}^{m_L-1} \frac{\partial z_j^L}{\partial a_k^{L-1}} \cdot \frac{\partial a_j^L}{\partial z_j^L} \cdot \frac{\partial Loss}{\partial a_j^L}$$

So, in the end:

$$\left[\frac{\partial \text{Loss}}{\partial w_{jk}^l} = a_k^{l-1} \cdot g'(z_j^l) \frac{\partial \text{Loss}}{\partial a_j^l} \right]$$

FOR US IS Δ_j^l

(WE USE GRADIENT DESCENT TO MINIMIZE THE LOSS)

$$\left[\begin{array}{l} \text{IS EQUAL TO } [2(y - a_j^l)] \\ \text{OR } \left[\sum_{j=\phi}^{M_{l+1}-1} w_{jk}^{l+1} \cdot g'(z_j^{l+1}) \frac{\partial \text{Loss}}{\partial a_j^{l+1}} \right] \end{array} \right]$$

ALL THE LOSS COMPOSE THE ∇_{Loss} !

GRADIENT DESCENT WITH BACKPROPAGATION

Let \tilde{X} the training set containing $[m]$ samples. So we have:

Loss \rightarrow Loss^(x) (loss for the sample $x \in \tilde{X}$) and

(WEIGHT UPDATE): $\underline{w}^l \leftarrow \underline{w}^l - d \times \frac{1}{m} \sum_x \underline{\Delta}^{l,x} (\underline{z}^{l-1}, x)^T$

ALL OF THOSE PARAMETERS APPEARS TO THE SAMPLE X

BIAS UPDATE: $\underline{w}_b^l \leftarrow \underline{w}_b^l - \frac{d}{m} \underline{\Delta}^l \cdot \underline{1}$

more simple $\left[\underline{w}^l \leftarrow \underline{w}^l - \frac{d}{m} \underline{\Delta}^l (\underline{z}^{l-1})^T \right]$

(COLUMN VECTOR OF ONES)

STOCHASTIC GRADIENT DESCENT

It's computationally faster split the training set T_s in subsets \tilde{X} containing m' samples each

1. Select a random subset \tilde{X} of m' samples from T_s
2. Train the network with \tilde{X} to estimate gradient (average over \tilde{X})
3. Repeat until there are not \tilde{X} s left

$\hookrightarrow (1+2+3) = \text{EPOCHS}$ of training \rightarrow repeat epochs until convergence.

NEURON SATURATION

Since we start with random weights, it is possible that weighted inputs to neurons are very big or very small at some point, close to the "stable" regions of threshold-like function.

This implies that in such case the derivative is ALMOST ZERO

Weights in hidden layers should be not initialized all to zero, in order to BREAK SYMMETRY

the modified error and the GRADIENT will be very small even if the output is not at all correct!

draw them from ~~an~~ INDEPENDENT GAUSSIAN distributions, with mean value = 0 and standard deviation = 1

But if a neuron has many inputs, the probability that it saturates is not negligible with this scheme \hookrightarrow a STANDARD DEVIATION = $\frac{1}{\sqrt{n}}$ is better ($n = \#$ INPUTS)

Cross-Entropy & Logistic Function

For an output layer with LOGISTIC NEURONS, we can design loss to eliminate saturation

\hookrightarrow Since $\Delta_j^L = \sigma'(z_j^L) \frac{\partial \text{loss}}{\partial a_j^L}$ We want $\frac{\partial \text{loss}}{\partial a_j^L} = \frac{a_j^L - y_j}{\sigma'(z_j^L) \sigma(z_j^L)(1 - \sigma(z_j^L))} = \frac{a_j^L - y_j}{a_j^L(1 - a_j^L)}$
 WE BUILT LOSS FOR HAVING THIS

$\rightarrow \int da_j^L \rightarrow \text{Loss}_j = -y \log(a_j^L) - (1 - y_j) \log(1 - a_j^L)$
 $\sigma'(z_j^L) = \sigma(z_j^L)(1 - \sigma(z_j^L))$

So $\left[\text{Loss} = -\sum_j y \log(a_j^L) + (1 - y_j) \log(1 - a_j^L) \right] \leftarrow$ CROSS-ENTROPY LOSS FUNCTION

It is always NON NEGATIVE, it is close to 0 when a_j^L is close to y (Recall that y and a^L are vectors of numbers between 0 and 1)

\hookrightarrow So, by construction, $\left[\Delta_j^L = a_j^L - y_j \right] \leftarrow$ WE BUILT "LOSS" FOR HAVING THIS

It is generally a BETTER solution than L_2 if the OUTPUT LAYER has LOGISTIC NEURONS!

THE SOFTMAX FUNCTION

If we want that the neurons of the output layers will return a PROBABILITY DISTRIBUTION that is not a binary outcome (logistic function) we have to use the SOFTMAX FUNCTION:

$\left[S(\vec{z})_i = \frac{e^{z_i}}{\sum_k e^{z_k}} \right] \leftarrow$ SOFTMAX ACTIVATION FUNCTION (each neuron represent an outcome and its activation represent the probability of that outcome)

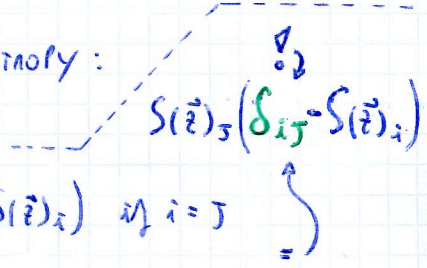
\hookrightarrow It takes in account the weighted input of ALL the neurons in the layer. It has values between 0 and 1, and the sum of the values along all the output layer is 1.

CROSS ENTROPY

We can use again, in a slightly different way, the CROSS-ENTROPY:

$\left[\text{Loss} = \sum_j y_j \log(a_j^L) \right]$

$(S(\vec{z})_i)$ is similar to $(\sigma(z_j^L))$, in fact: $\left(\frac{\partial S(\vec{z})}{\partial z_j} \right)_i = \begin{cases} S(\vec{z})_j (1 - S(\vec{z})_i) & \text{if } i = j \\ -S(\vec{z})_i S(\vec{z})_j & \text{otherwise} \end{cases}$



Knowing that, and remembering that $a_i^L = S(\tilde{z}^L)_i$, we can compute the gradients:

$$\frac{\partial \text{loss}}{\partial w_{jk}^L} = \sum_i \frac{\partial \text{loss}}{\partial a_i^L} \frac{\partial a_i^L}{\partial w_{jk}^L} = \sum_i y_i \frac{\partial a_i^L}{\partial z_j^L} \frac{\partial z_j^L}{\partial w_{jk}^L} = \sum_i y_i \frac{\partial a_i^L}{\partial z_j^L} a_k^{L-1} = \sum_i y_i \frac{\partial a_i^L}{\partial z_j^L} (s_{ij}^L a_i^L) a_k^{L-1}$$

[But \vec{y} is $\emptyset \forall i \neq j$ (null for each component representing a wrong answer)] \checkmark s_{ij}^L

$$\frac{\partial \text{loss}}{\partial w_{jk}^L} = \frac{y_j}{s_j^L} a_j^L (1 - a_j^L) a_k^{L-1} = a_k^{L-1} \underbrace{(y_j - a_j^L)}_{\Delta_j^L} \quad \text{OK NO SATURATION (AS WITH LOGISTIC FUNCTION)}$$

REGULARISATION IN ANN

It's helpful to use REGULARISATION to limit OVERFITTING (AS SEEN IN 4.2 LINEAR MODELS) but this time we can use COMBINATIONS of L1+L2 REGULARIZATION \rightarrow ELASTIC NET OF REGULARIZATION. PL(42)

This can also be done by acting directly on the weights \rightarrow DRAGOUT REGULARIZATION

(As with all other methods, INCREASING THE SIZE of the training set, possibly artificially,) reduces overfitting.)

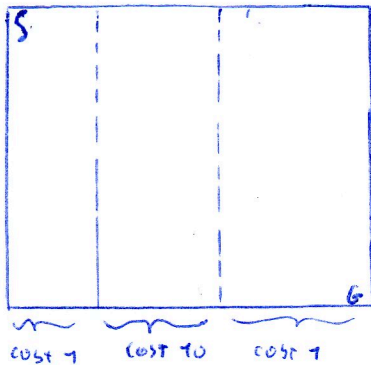
ARTIN EXERCISES

- EXERCISES 49-51 - EX-1: A*

Robot moving in 12x10 tiles, with random walls. MOVES: ↑, ↓, →, ←,

START: (1,1) GOAL: (12,10)

COST of PASSING = $\begin{cases} \text{cols } 4,5,6,7 \Rightarrow \textcircled{10} \\ \text{cols } 1,2,3,8,9,10,11,12 \Rightarrow \textcircled{1} \end{cases}$



Q1 Assuming that exists a path to the goal, can the robot forced to move over more than 4 of the tiles that have cost 10?

Yes, if there is not a full row free from column 4 to column 8 you cannot have a straight line between the "cost 10" zone

Q2 Give an admissible heuristic for the robot to find a cost-optimal path in this problem, not trivial and take in account the fact that cols 4-7 cost is 10.

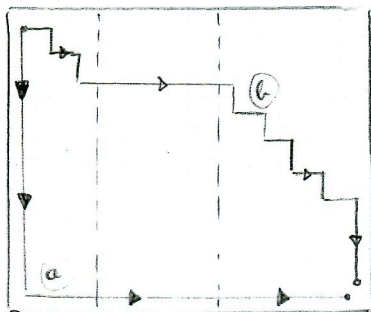
$h(x) = \text{manhattan}(x, \text{GOAL}) = \Delta x + \Delta y \rightarrow$ have to take account of the cost

$\hookrightarrow h(x) = \text{man}(x, \text{GOAL}) + \text{col}(x)$ with

$$\text{col}(x) = \begin{cases} (10-1) \cdot 4 & \text{if } x < 4 \\ \emptyset & \text{if } x > 7 \\ (10-1) \cdot (7-x) & \text{if } x \in [4, 7] \end{cases}$$

CAUSE 1 IS COUNTED ALREADY IN MANHATTAN

Is it admissible?



Yes, because it never overestimate the cost of a x .

$$m(x, y) + \text{col}(x)$$

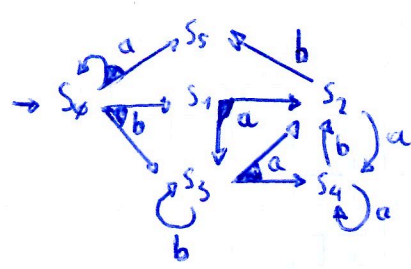
if $x < 4$, 4

if $x > 7$ 0

if $x \in [4, 7]$ $(7-x) \cdot 9$

- Exam '17-'18 - Ex2: Non-determinism

Q3 Does there exist a winning strategy to go from s_0 to s_5 in the following hypergraph?



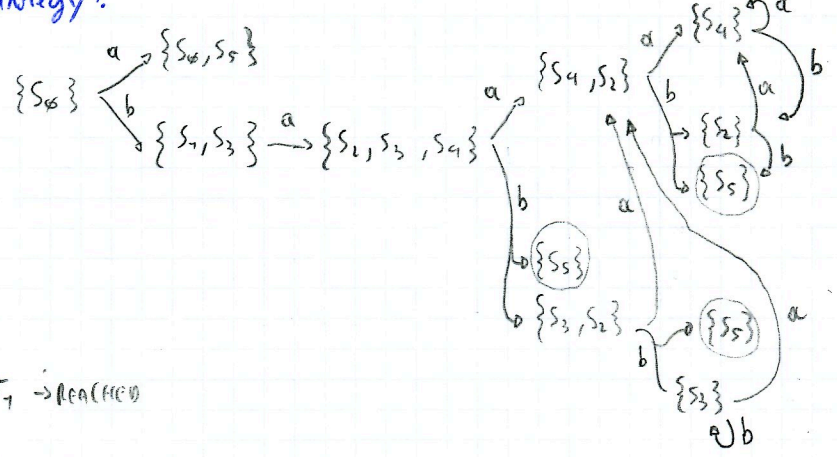
$W = s_5$
 $R = [F, F, F, F, F, T]$
 ↓ for s_5
 $R[s_0] = F$ $R[s_2] = T \Rightarrow \text{strat}[s_2] = b \Rightarrow W = s_2$
 ↓ for s_2
 $R[s_1] = F$ $R[s_4] = T \Rightarrow \text{strat}[s_4] = b \Rightarrow W = s_4$
 ↓ for s_4
 $R[s_3] = T \Rightarrow \text{strat}[s_3] = a \Rightarrow W = s_3$
 ↓ for s_3
 $R[s_1] = T \Rightarrow \text{strat}[s_1] = a$ $R[s_0] = T \Rightarrow \text{strat}[s_0] = b$

WINNING strategy:

- $s_0 \rightarrow b$
- $s_1 \rightarrow a$
- $s_2 \rightarrow b$
- $s_3 \rightarrow a$
- $s_4 \rightarrow b$
- $s_5 \perp$

Q4 Suppose now the system is partially observable: $(O(s_0) = O(s_5) = \emptyset)$ and $(O(s_1) = O(s_2) = O(s_3) = O(s_4) = \{a\})$. Does exist a winning strategy?

Respective Belief Hypergraph:



WINNING strategy (example):

$\sigma_1 \rightarrow (b) \rightarrow \sigma_2 \rightarrow (a) \rightarrow \sigma_2 \rightarrow (b) \rightarrow \sigma_1 \rightarrow \text{REACHED}$
 $\sigma_2 \rightarrow (b) \rightarrow \sigma_1 \rightarrow \text{REACHED}$
 $\sigma_2 \rightarrow (a) \rightarrow \sigma_2 \rightarrow (b) \rightarrow \sigma_1 \rightarrow \text{REACHED}$
 $\sigma_2 \rightarrow (b) \rightarrow \sigma_1 \rightarrow \text{REACHED}$

but an imconclusive one is $\left[\sigma_1 \rightarrow (b) \rightarrow \sigma_2 \rightarrow (a) \rightarrow \sigma_2 \rightarrow (a) \rightarrow \sigma_2 \rightarrow (a) \rightarrow \sigma_2 \rightarrow (b) \rightarrow \sigma_2 \rightarrow (b) \rightarrow \sigma_1 \right]$

- Exam 47-18 - Ex 3: Markov Processes

1	2	3
1	-0.07	-1
2	-0.07	+2
3	0.07	-0.07

FINAL

10% of going LEFT instead of going in the desired direction
(if goes into wall stays there)

Q5: Compute the first 3 iterations of the VALUE ALGORITHM for the MDP
Assume $\gamma = 1$. Detail the computation for state (2,2)

$$\begin{aligned}
 \left(\begin{array}{l} \uparrow \Rightarrow \uparrow 90\% + \leftarrow 10\% \\ \leftarrow \Rightarrow \leftarrow 90\% + \downarrow 10\% \\ \downarrow \Rightarrow \downarrow 90\% + \rightarrow 10\% \\ \rightarrow \Rightarrow \rightarrow 90\% + \uparrow 10\% \end{array} \right) & U_{i+1}(s) = R(s) + \gamma \max_{a \in A} \sum_{s' \in S} P(s', a, s) U(s') \\
 U_{\emptyset}(s_i) &= \emptyset \quad \forall s_i
 \end{aligned}$$

$$U_1(2,2) = -0.07 + \left[\uparrow 0.9(\emptyset) + 0.1(\emptyset), \leftarrow 0.9(\emptyset) + 0.1(\emptyset), \downarrow 0.9(\emptyset) + 0.1(\emptyset), \rightarrow 0.9(\emptyset) + 0.1(\emptyset) \right] = -0.07$$

$$U_1(s_i) = R(s_i), \forall s_i$$

$$U_2(2,2) = -0.07 + \left[\uparrow 0.9(-0.07) + 0.1(-0.07), \leftarrow 0.9(2) + 0.1(-0.07), \downarrow 0.9(-0.07) + 0.1(-0.07), \rightarrow 0.9(-0.07) + 0.1(-0.07) \right] = -1.723$$

$$U_2(1,1) = -0.14 \quad \text{strat } [1,1] = \leftarrow ; \quad U_2(2,1) = -0.14 \quad \text{strat } [1,2] = \uparrow$$

$$U_2(2,3) = -0.14 \quad \text{strat } [2,3] = \downarrow ; \quad U_2(3,2) = 2 ; \quad U_2(3,3) = -1$$

$$\hookrightarrow \text{STRAT: } (1,1) \leftarrow [-0.14] / (1,2) \uparrow [-0.14] / (2,2) \leftarrow [-1.723] / (3,2) \uparrow [2] / (2,3) \downarrow [-0.14] / (3,3) \uparrow [-1]$$

Q6 Suppose STRAT: $\leftarrow \forall s_i, P(\text{SQUEAK IN } (j,i)) = 10 \cdot (j+2)\%$

$$X = \begin{bmatrix} (3,1) \\ (3,2) \\ (2,2) \\ (1,2) \\ (2,3) \\ (1,3) \end{bmatrix} \quad P = \begin{bmatrix} 1/4 \\ 1/4 \\ 1/4 \\ 1/4 \\ \emptyset \\ \emptyset \end{bmatrix}$$

We observe (Squeak, Squeak, Not Squeak) What is the most probable position of the robot.

	(3,1)	(3,2)	(3,3)	(1,2)	(2,3)	(1,3)
(3,1)	0.1	0.9	\emptyset	\emptyset	\emptyset	\emptyset
(3,2)	\emptyset	0.9	0.1	\emptyset	\emptyset	\emptyset
(2,2)	\emptyset	\emptyset	\emptyset	0.1	0.9	\emptyset
(1,2)	\emptyset	\emptyset	\emptyset	0.1	\emptyset	0.9
(2,3)	\emptyset	\emptyset	\emptyset	\emptyset	1	\emptyset
(1,3)	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset	1

$$P_{\text{SQUEAK}} = \begin{bmatrix} 40\% & \emptyset & \emptyset & \emptyset & \emptyset & \emptyset \\ \emptyset & 50\% & \emptyset & \emptyset & \emptyset & \emptyset \\ \emptyset & \emptyset & 40\% & \emptyset & \emptyset & \emptyset \\ \emptyset & \emptyset & \emptyset & 30\% & \emptyset & \emptyset \\ \emptyset & \emptyset & \emptyset & \emptyset & \emptyset 90\% & \emptyset \\ \emptyset & \emptyset & \emptyset & \emptyset & \emptyset & \emptyset 10\% \end{bmatrix}$$

$$Q_{1:t} = (S, S, NS)_2$$

$$l_{1:1} = Q_S T^T l_\varnothing = \begin{bmatrix} 0.07 \\ 0.225 \\ 0.07 \\ 0.0150 \\ 0 \\ 0 \end{bmatrix}$$

$$l_{1:2} = Q_S T^T l_{1:1} = \begin{bmatrix} 0.0004 \\ 0.1058 \\ 0.0090 \\ 0.0008 \\ 0 \\ 0 \end{bmatrix}$$

$$l_{1:3} = Q_{NS} T^T l_{1:2} = \begin{bmatrix} 0.0000 \\ 0.0478 \\ 0.0063 \\ 0.0007 \\ 0.0087 \\ 6.0707 \end{bmatrix}$$

$$F_{1:1} = \frac{l_{1:1}}{\text{sum}(l_{1:1})} = \begin{bmatrix} 3.8462\% \\ 86.5385\% \\ 3.8462\% \\ 5.7692\% \\ 0\% \\ 0\% \end{bmatrix}$$

$$F_{1:2} = \begin{bmatrix} 0.3451\% \\ 91.2425\% \\ 7.7653\% \\ 0.6471\% \\ 0\% \\ 0\% \end{bmatrix}$$

$$F_{1:3} = \begin{bmatrix} 0.0377\% \\ 75.1132\% \\ 9.9774\% \\ 1.0732\% \\ -12.7371\% \\ 1.0614\% \end{bmatrix}$$

(WITH INTERS)

$$\begin{bmatrix} \approx 75\% \text{ that robot in im (3,2)} \\ \approx 13\% \text{ im (2,3)} \\ \approx 10\% \text{ im (2,2)} \end{bmatrix}$$

- EXAM 17-18 - Ex 4: SUPERVISED LEARNING

x_1	0	1	0	1	1	0	1	1	0	0
x_2	1	1	0	1	0	0	0	1	1	0
x_3	1	1	1	0	0	1	1	1	0	0
$f(x)$	1	1	0	1	1	0	1	1	0	0

Q3 What is the IG for x_1 and x_2 ? Which one should be chosen as the root for a decision tree for approximating f ?

$$IG(A) = H_B\left(\frac{p}{p+m}\right) - R(A) \quad R(A) = \sum_{k=1}^d \frac{p_k+m_k}{p+m} H_B\left(\frac{p_k}{p_k+m_k}\right) \quad H_B(q) = -[q \log_2(q) + (1-q) \log_2(1-q)]$$

x_1) $\begin{matrix} \text{0} & \text{1} & \text{0} & \text{1} & \text{1} & \text{0} & \text{1} & \text{1} & \text{0} & \text{0} \\ \text{1} & \text{1} & \text{0} & \text{1} & \text{0} & \text{0} & \text{1} & \text{1} & \text{0} & \text{0} \end{matrix} \rightarrow \begin{matrix} p_0+m_0=5 & p_0=1 \\ p_1+m_1=5 & p_1=5 \end{matrix} \quad H_B\left(\frac{1}{5}\right) = H_B\left(\frac{3}{5}\right) = -\left[\frac{3}{5} \log_2\left(\frac{3}{5}\right) + \frac{2}{5} \log_2\left(\frac{2}{5}\right)\right] = 0.97$

$$H_{B_0}\left(\frac{1}{5}\right) = -\left[1 \cdot \frac{1}{5} \log_2\left(\frac{1}{5}\right) + \frac{4}{5} \log_2\left(\frac{4}{5}\right)\right] = 0.72 \quad H_{B_1}\left(\frac{5}{5}\right) = 0$$

$$R(x_1) = \frac{1}{2} \cdot 0.72 + \frac{1}{2} \cdot 0 = 0.36 \quad IG(x_1) = 0.97 - 0.36 = 0.61$$

x_2) $\begin{matrix} \text{1} & \text{1} & \text{0} & \text{1} & \text{0} & \text{0} & \text{1} & \text{1} & \text{0} & \text{0} \\ \text{1} & \text{1} & \text{0} & \text{1} & \text{1} & \text{0} & \text{1} & \text{1} & \text{0} & \text{0} \end{matrix} \rightarrow \begin{matrix} p_0+m_0=5 & p_0=2 \\ p_1+m_1=5 & p_1=4 \end{matrix}$

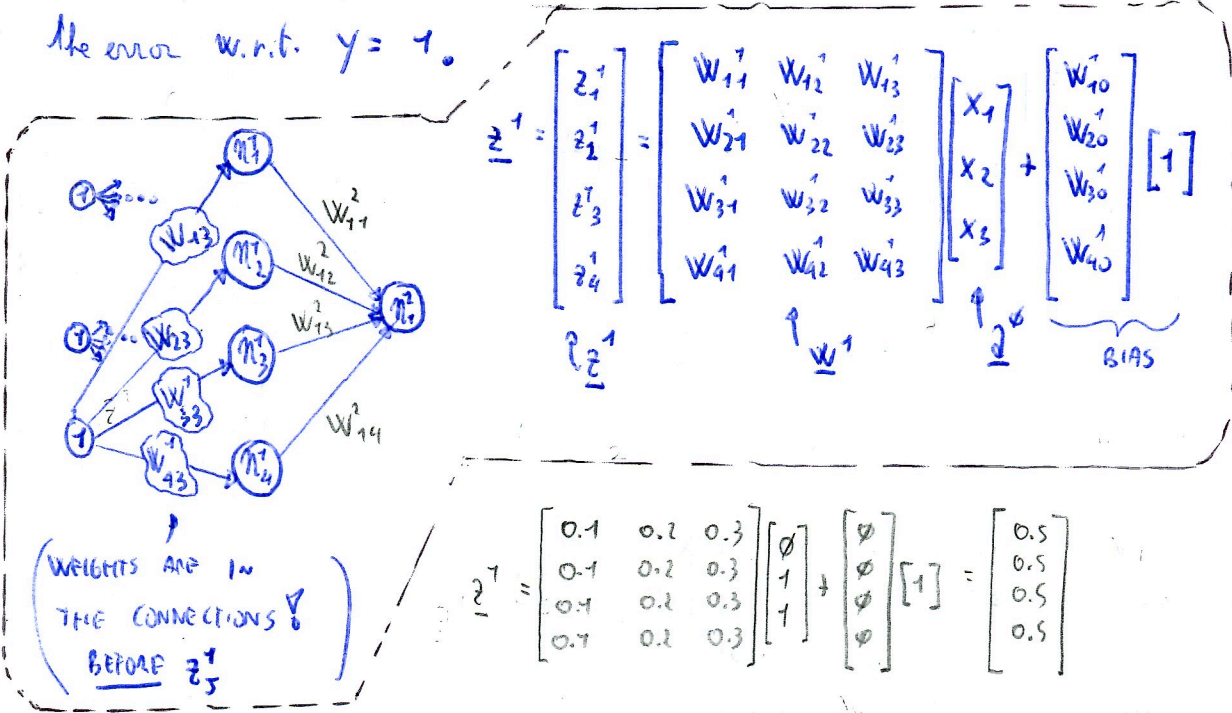
$$H_{B_0}\left(\frac{2}{5}\right) = -\left[\frac{2}{5} \log_2\left(\frac{2}{5}\right) + \frac{3}{5} \log_2\left(\frac{3}{5}\right)\right] = 0.97 \quad H_{B_1}\left(\frac{4}{5}\right) = 0.72$$

$$R(x_2) = \frac{1}{2} \cdot 0.97 + \frac{1}{2} \cdot 0.72 = 0.845 \quad IG(x_2) = 0.125$$

$[IG(x_1) > IG(x_2)] \Rightarrow$ ROOT NODE will test x_1 (more information on x_1)

Q8 Now we want to approx an ANN with a single hidden layer of 4 neurons. We use cross entropy + logistic function. Biases initially 0.

$W_{JK}^l = (k+1-l) \cdot 0.1$. Compute the prediction for $(0, 1, 1)$ and the error w.r.t. $y = 1$.



$$z^1 = \begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.1 & 0.2 & 0.3 \\ 0.1 & 0.2 & 0.3 \\ 0.7 & 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.5 \\ 0.5 \\ 0.5 \end{bmatrix}$$

$$W^2 = [W_{11}^2 \ W_{12}^2 \ W_{13}^2 \ W_{14}^2] = [0 \ 0.1 \ 0.2 \ 0.3]$$

$$f(0.5) = \frac{1}{1 + e^{-(0.5)}} = 0.622 \quad \Rightarrow \quad a^1 = \begin{bmatrix} 0.622 \\ 0.622 \\ 0.622 \\ 0.622 \end{bmatrix}$$

REMEMBER TO PASS THE z THROUGH THE SIGMOID FUNCTION TO GET a !

$$z^2 = W^2 \cdot a^1 + [W_{10}^2] \cdot [1] = 0.622(0.1 + 0.2 + 0.3) - 0.1 = 0.3732 - 0.1 = 0.2732$$

$(0 - 1) \cdot 0.1 = -0.1$

$$a^2 = f(z^2) = \frac{1}{1 + e^{-(0.2732)}} = 0.5679 \quad \Delta^2 = a^2 - y = 0.5679 - 1 = -0.4321$$

$$(\Delta^l = (W^l)^T \Delta^l \cdot g'(z^{l-1}); \quad g'(z^{l-1}) = g(z^{l-1}) (1 - g(z^{l-1})))$$

$$\Delta^1 = \begin{bmatrix} 0 \\ 0.1 \\ 0.1 \\ 0.3 \end{bmatrix} [-0.4321] \cdot \begin{bmatrix} 0.25 \\ 0.25 \\ 0.25 \\ 0.25 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.025 \\ 0.025 \\ 0.075 \end{bmatrix} [-0.4321] = \begin{bmatrix} 0 \\ -0.0108 \\ -0.0216 \\ -0.6324 \end{bmatrix} = \underline{\Delta^1}$$

