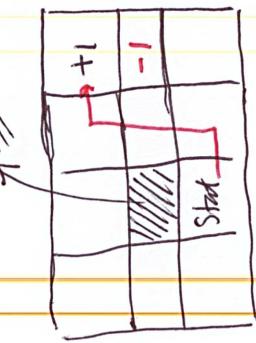


SYSTEM

Sensors — States — actions

Prob (action) < 1

Example

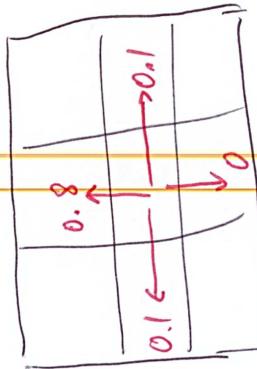


Objectives
Reach +1 from START
w/o passing through
-1

Answer: shortest path that avoids -1
This is a deterministic model

Move from deterministic to a non-deterministic (probabilistic)
model: for each cell assign a probability
to move N/E/S/W

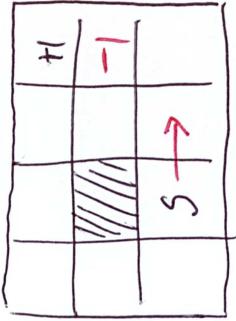
Now $P(N|Start)$ is max
and, since we chose to
provide more by prob.,
the more w'll be North,
where there is a wall!



So, it is NOT a
good option!

Let us consider the other two options.

(1) More East



P-20.1

(12) More West



$$p \geq 0.1$$

\Rightarrow Prob : 0

$$\text{Prob} : 0.1 \times 0.8 \times 0.8 \times 0.1 \times 0.1 \approx 0.1 \rightarrow \text{hand } +1$$

Dr. Balsoreo,

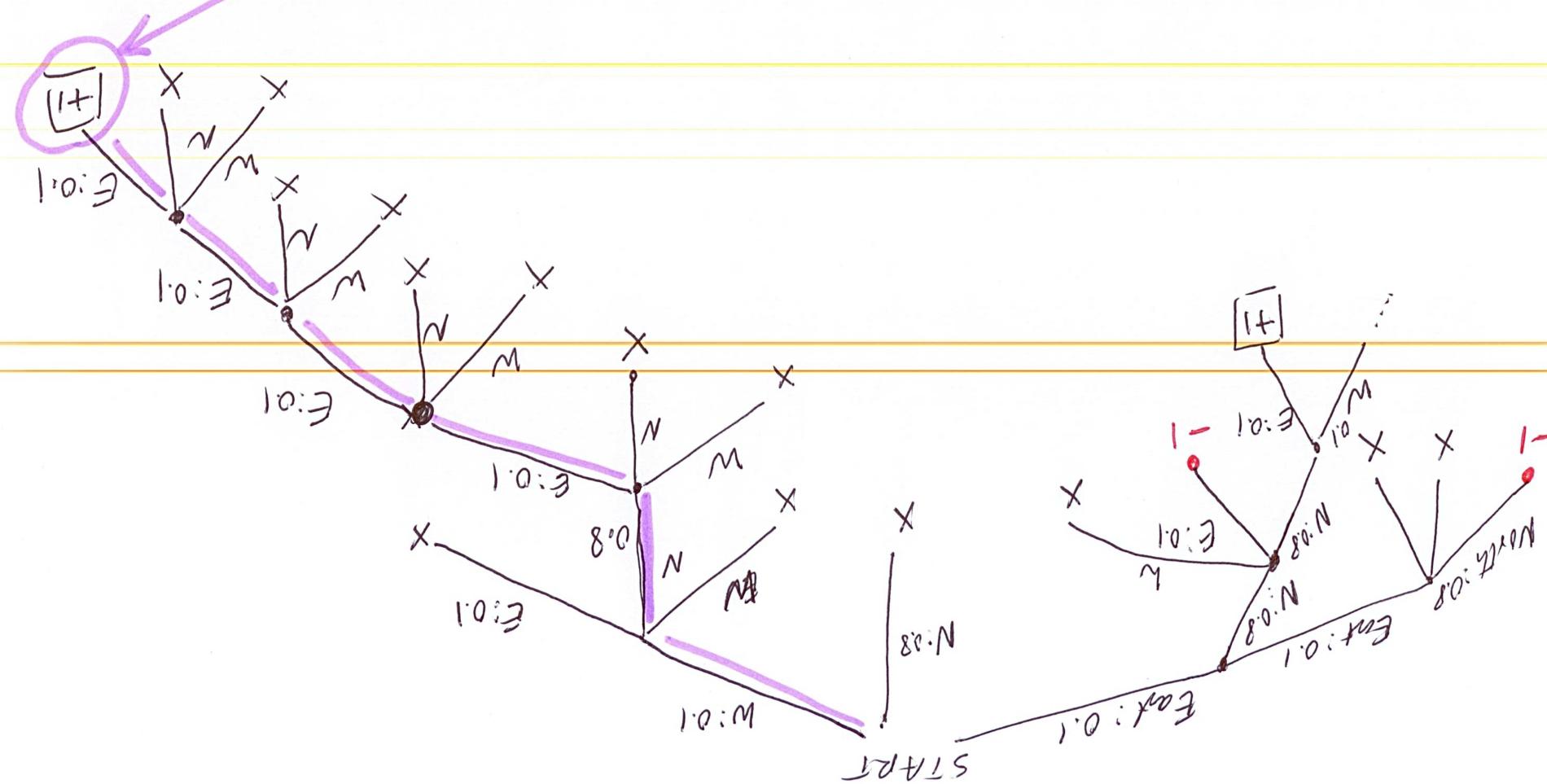
We can illustrate all the moves / paths with their prototypes !

MDP-3

But what is what?

(Highest Utility)

To reach +1
longer path



In the Deterministic Transition Model

We use the shortest paths (A^* -alg., Dijkstra)
 Utility is a function from the current state
 (a node in the graph of states) to the goal state.

We want the utility to be inverse related to the goal state:
smallest distance \equiv highest utility

$$U(\text{node}) = \frac{1}{\text{dist}(\text{node}, \text{goal})}$$

The DT Model assumes a perfect world

From utility to Policy: mapping from states to actions
 Policy: States \rightarrow Actions / π is used of few to denote Policy
Objective: Find "BEST policy"

Markov Decision Problem: Given a stochastic environment with transition model known, compute the optimal policy

Markov Property

MDP-5

s_i, s_{i-1}, \dots is a sequence of states -

$$\text{MP: } P(s_i | s_{i-1}, \dots, s_0) = P(s_i | s_{i-1})$$

NOT EVERY DECISION PROBLEM IS A MDP.

Now, let us formalize the MDP assuming that we define the following

$P_a(i, j)$: prob. to reach state j from state i by action a
 $U(j)$: utility of node/state j

Then for a collection of nodes $J \subseteq \mathcal{T}$, and actions $a \in A$

$$J_{\text{opt}}(i) = \arg \max_{a \in A} \sum_j P_a(i, j) U(j)$$

In general we do NOT know U !

To compute U we need to look at the tree of states. We want to estimate how good a state is.



Obviously, utility of state is as good as the utilities of its successor states are.

$$s_0 \rightarrow s_1 \rightarrow \dots \rightarrow s_n$$

$$o \rightarrow a \rightarrow 2 \rightarrow \dots \rightarrow n$$

$$U(s) = f(U(s_1), \dots, U(s_n)) \text{ if } r$$

Separable utility: utility can be decomposed into two components:

- local component : REWARD
- non-local / successor component : utility of successors

$$U([0, 1, \dots, n]) = R(0) + U([1, \dots, n])$$

In general,

$$U([i, i+1, \dots, n]) = R(i) + U([i+1, \dots, n])$$

or,
$$\boxed{U(i) = R(i) + \max_a \sum_j p_{a(i,j)} U(j)} \quad (1)$$

Note $R(i)$, the local component of $U(i)$, does not depend on the successor states. In a way it captures what we mean when we say "So far so good", or "easy does it".

Equation (1) captures the OPTIMALITY PRINCIPLE from Dynamic Programming.

"An optimal policy has the property that whatever the initial state and initial actions are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first action."

For an n -step decision problem:

Naive solution: $O(|A|^n)$ A : set of actions.

DP solution $O(n|A||S|)$ S : set of states

What is the "correct" value of n ?

Computing the utility function

Idea for iterative computation: $U_t(i)$ util'g of tree t

$$U_{t+1}(i) = R(i) + \max_{a \in A} \sum_j P(a, j) U_t(j)$$

Then define the optimal utility

$$U_{opt} = \lim_{t \rightarrow \infty} U_t$$

U_{t+1}

Operationally,

$$U_{opt} = \begin{cases} U_t & \text{if } U_{t+1} > U_t \\ U_{opt} = U_{t+1} & \text{else if } U_{t+1} < \text{threshold} \\ \text{STOP} & \end{cases}$$

Formally, we have

Given $(p^{(i,j)})_{i,j}$; a transition matrix

R : a reward function; ϵ : threshold
we calculate the utility as follows

Initialize U, U_1 to R

Repeat $U = U'$

for $i = 1 : n$

$$U(i) = R(i) + \max_a \sum_j p^{(i,j)} U(j)$$

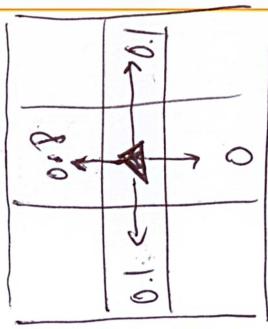
end

until $|U - U'| < \epsilon$

Return U

Value Iteration Alg

Return to example to compute U



	$u=10$	
$u=5$	$u=1$	$u=2$
	$u=1$	

u, R

$Q \in \{N, W, E, S\}$

$$a : w \leftarrow 0.1 \times 1 + \boxed{0.8 \times 5} + 0.1 \times 10 = 5.1$$

$$a : N \uparrow 0.1 \times 5 + \boxed{0.8 \times 10} + 0.1 (-8) = 7.7$$

$$a : E \rightarrow 0.1 \times 1 + 0.1 \times 10 + \boxed{0.8 (-8)} = -5.3$$

$$a : S \downarrow 0.1 (-8) + 0.1 \times 5 + \boxed{(0.8)(11)} = 0.5$$

$$U(\text{center}) = R(\text{center}) + \max \{ 5.1(w), 7.7(N), (-5.3)(E), 0.5(S) \}$$

$$= 1 + 7.7(N) = 8.7(N)$$

Utility for center is 8.7 (going up).

We can now compute the optimal policy

$$\pi_{opt}(i) = \arg \max_a \sum_j P_{a(i,j)} \pi_{opt}(j)$$

This is the optimal solution for MDP.

ISSUES: How to measure the convergence?
Use some measure such as Root Mean Square Error

$$RMSE = \frac{1}{|S|} \left(\sum_{i=1}^{|S|} (u(i) - u(i))^2 \right)^{\frac{1}{2}}$$

Stop when $RMSE(u, u') < \epsilon$
Don't use ϵ :

0.81	0.86	0.91	+1
0.76	0.77	0.66	-1
0.70	0.65	0.61	0.68

Calculate utilities: u_{opt}

+1			

execute actions:

π_{opt}

+1			
-1			

original env.

+1			
-1			