



Fakultät für Informatik
Lehrstuhl für Echtzeitsysteme und Robotik

Get Me Out Of Here: Determining Optimal Policies

Marcel Bruckner

Seminar *Cyber-Physical Systems* WS 2017/18

Advisor: Christian Pek

Supervisor: Prof. Dr.-Ing. Matthias Althoff

Submission: 01. January 2017

Get Me Out Of Here: Determining Optimal Policies

Marcel Bruckner
Technische Universität München
Email: mbruckner94@gmail.com

Abstract—This paper gives a brief overview on the field of dynamic programming and it's applicability in robot motion planning. Exemplary, some use cases in other papers will be shown and an algorithm to find optimal policies in a maze will be discussed.

I. INTRODUCTION

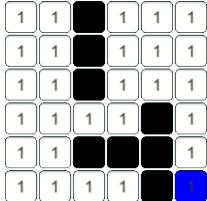
Remember the game Labyrinth, in which the game represents a maze where the field is build up by fixed and moving pieces showing walls and passages. The goal of the game is to rearrange the maze by moving rows of not fixed tiles and enable ones game character to find treasures and return them to goal tiles.

To reach the goal in this game, a way for calculating the shortest path to exit the maze has to be found. This can be achieved with classical motion planning algorithms, like Dijkstras Algorithm or Breadth-First Search, but also with **Dynamic Programming**, which will be used in this Paper.

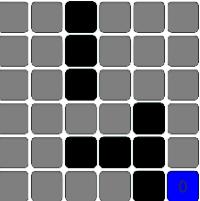
Therefore, let a grid map $M = \mathbb{N}^{n \times m}$ containing static walls $M(i, j) = \infty$ represent the maze and let there be a set of exits X_{goal} in the labyrinth which the player wants to reach.

The goal of this paper is to point out ways to determine optimal policies $P = \mathbb{N}^{n \times m} \rightarrow A$ to exit the labyrinth. To achieve this, a way to calculate the value $v(i, j)$ of the current position is introduced from what the actual best action can be derived.

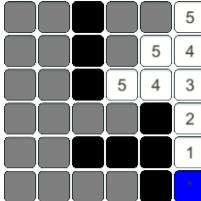
Later on, the algorithm will look like in the following example:



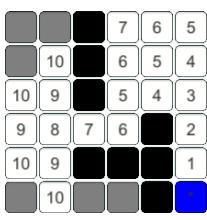
The original grid showing the cost of each tile, walls (black) and goals (blue).



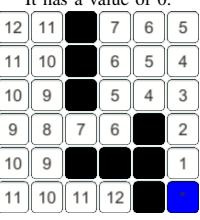
Initialization of the algorithm with ∞ cost. Only the goal (blue) is set and is the starting point. It has a value of 0.



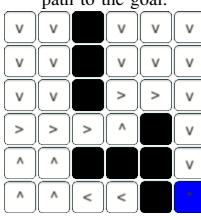
After a few iterations the values of some tiles are set. The value is equal to the length of the shortest path to the goal.



A few more iterations. The algorithm is reaching tiles further away from the goal.



All values are calculated. Every tile is holding its value.



The final optimal policy showing the best action for every position derived from the values of the tiles.

II. FUNDAMENTALS

A. What is Dynamic Programming

Dynamic Programming is a method used in mathematical optimization and computer programming that transforms complex problems into sequences of simpler subproblems. Each of those subproblems will be solved just once and their solutions will be stored, so when reoccurring, the previously computed solution will be looked up to save computation time.

For solving a problem with Dynamic Programming, the problem gets broken down and by recursively finding the optimal solution for the subproblems a optimal solution for the whole problem can be found.

It provides a general framework for analyzing many problem types, like the planning of industrial production lines, scheduling of patients at a medical clinic, long term investment programs and, of course, robot motion planning.

//TODO: Mehr ausbauen Review ARSLAN DISSERT P.13
Underlying of these problems is always a sequence of decisions that have to be made within a finite number of stages. E.g. a robot wants to move from one position to another by a finite number of steps and investment programs do have to calculate their wealth at every step in time.

"A sequential decision making problem with finite stages is a tupel $\langle X, U, f, g, G, x_{\text{init}} \rangle$ where

- X is the set of all possible states of the system, i.e., the state space;
- U is the set of controls, i.e., the control space;
- $f : X \times U \mapsto X$ is a transition function, a mapping specifying the next state $f(x, u)$ if control u is executed when the system is in state x ;
- $g : X \times U \mapsto \mathbb{R}$ is a cost function that gives a real value $g(x, u)$ obtained if control u is executed when the system is in the state x ;
- $G : X \mapsto \mathbb{R}$ is a cost function that gives a real value $G(u)$ obtained if the final state is x ;
- $x_{\text{init}} \in X$ is the initial state of the system." [10]

B. Richard Bellman

Dynamic Programming was introduced by the american mathematician Richard Bellman who was working on multistage decision problems from 1949 until 1955. After he received his Ph.D. from Princeton at the age of 25, he started his research at RAND Corporation (Research and Development), where his secretary really had a hatred of the term *mathematical research*. So he started his work with finding a name under which he could hide his projects.

He was interested in planning, decision making and thinking, but none of these were good terms, so he decided on the word programming, and because of his love for physics and his ideas in multistage and time varying problems, he came to the conclusion that dynamic is the best adjective to accompany it. Following the name Dynamic Programming was born and no one could imagine what he was doing which made it perfect for hiding his mathematical research.

C. The Principle of Optimality

For starting a more precise discussion, let define the term optimal policy:

- A **policy** is a sequence of decisions.
- An **optimal policy** is a sequence of decisions that has the best outcome w.r.t. a predefined criterion.

When being asked the task for finding the optimal policy for a given problem, the classical approach would be to compute all possible policies and maximize the return to determine which of these is the optimal one. In practice, this is often not practical due to the fact, that even for a low number of stages and choices the dimension of the resulting maximization problem will explode.

But Bellman saw that the knowledge of the complete sequence of decisions isn't even required to solve these multistage decision problems. He stated that it would be much better to find a general prescription which gives the best decision at any time and depending only on the current state of the system.

"If at any particular time we know what to do, it is never necessary to know the decisions required at subsequent times." [3]

1) The fundamental approach:

From the previous stated point of view, Bellman went over to define an optimal policy and described it as one "determining the decision required at each time in terms of the current state of the system. Following this line of thought, "[3] he was able to formulate his

Principle of Optimality

"An optimal policy has the property, that whatever the initial state and initial decisions are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decisions." [3]

and thereof derive the fundamental mathematical equations on which Dynamic Programming is based on.

2) Mathematical formulation:

But what does this mean for the functional equations needed for actual problem solving based on the principle of optimality?

Imagine the simplest case, where a process of a system is described at any time by an M-dimensional vector

$$p = (p_1, p_2, \dots, p_M) \in X \quad (1)$$

Let $U = \{T_k\}$ be a set of transformations with the property

$$p \in X \Rightarrow T_k(p) \in X, \quad \forall k \quad (2)$$

Now thinking of an N-stage process, the goal is to maximize a scalar function $R(p)$ giving the return or the value of the final state. This function will be called the N-stage return and a policy consists of a selection of N transformations

$$P = (T_1, T_2, \dots, T_N) \quad (3)$$

yielding successively the states

$$\begin{aligned} p_1 &= T_1(p), \\ p_2 &= T_2(p_1), \\ &\vdots \\ p_N &= T_N(p_{N-1}) \end{aligned} \quad [3] \quad (4)$$

As you can see the maximum value of the return function $R(p_N)$, regulated by the optimal policy, will only depend on the initial vector p and the number of stages N.

One can now give a direct formula for the N-stage return that we get only using the optimal policy and by starting from the initial state p as

$$h_N(p) = \max_P R(p_N) \quad (5)$$

To get a functional equation for $h_N(p)$, the principle of optimality is used. By making our first decision, some transformation T_k will be made and we come to a new state $T_k(p)$. From the new state $T_k(p)$, there are now left $(N - 1)$ stages and the maximum return of these is, as defined above, $f_{N-1}(T_k(p))$.

As you might see, k and thereof T_k have to be chosen so as to maximize the return, which finally results in the basic functional equation

$$h_N(p) = \max_k h_{N-1}(T_k(p)), \quad N = 2, 3, \dots \quad (6)$$

Resulting from this equation it is now clear, that any optimal policy does not have to be unique, but will yield $h_N(p)$. On the other hand, if looking at the sequence $\{h_N(p)\}$, one can reversely determine all optimal policies.

3) Which characteristics does a problem need to be solvable by Dynamic Programming:

a) Overlapping subproblems:

A problem is said to have overlapping subproblems if the space of subproblems is small and thus, if broken down, solutions for these are reused many times or a recursive algorithm solves the same subproblems over and over, instead of generating new subproblems.

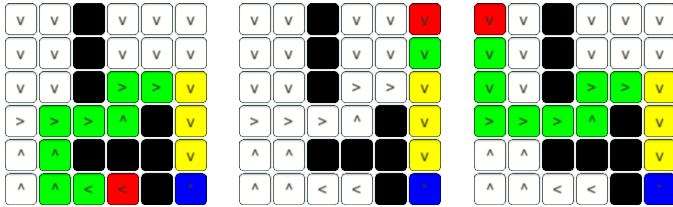
Remember the game Labyrinth, where the goal was to find treasures and return to a goal tile.

If there is a treasure you want to discover and it lays down a narrow passage, any algorithm calculating the shortest path from any arbitrary position in the maze will guide you through

the passage.

So if you would calculate the shortest path from your current position and move one tile, you'd have to recalculate the path again, which leads you to also recalculate the way down the passage.

So for any position, the passage would be a subproblem that overlaps in every calculation.



An example of overlapping subproblems. In all three cases the yellow passage will be recalculated no matter from where you started (red).

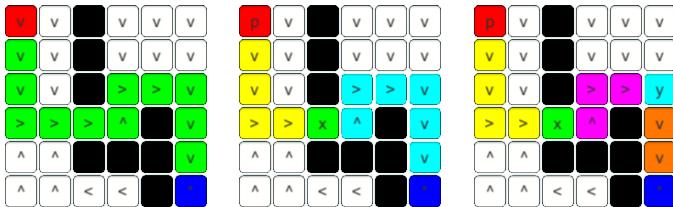
b) Optimal substructure:

A problem has optimal substructure if the solution for the whole problem can be found by a combination of solutions for its subproblems. Normally, greedy algorithms are used to solve problems with this property because they lead to an optimal overall solution in this case.

But also the principle of optimality is based in this idea.

To come back to the labyrinth, for any given position p in the labyrinth a shortest path s avoiding the walls and reaching a goal can be calculated. If s is really the shortest path, it can be split into subpaths s_1 from p to an point x on the shortest path and s_2 from x to the goal, such that s_1, s_2 are indeed the shortest paths between the three positions.

From that, a recursive way for finding shortest paths can easily be formulated, which is what is applied in Dynamic Programming.



The overall shortest path s (green) from start (red) to goal (blue).

A first split of s into s_1 (yellow) and s_2 (cyan).

A second split of s_2 into $s_{2,1}$ (magenta) and $s_{2,2}$ (orange).

An example of optimal substructure. The shortest path s (green) can be split down into subproblems recursively until only the shortest path between two neighbored tiles has to be calculated. Based on these optimal solutions for the subproblems, the whole problem can be solved.

4) How to apply Dynamic Programming in computer science:

a) Memoization:

An important concept when applying Dynamic Programming in computer science is memoization. As a result of the optimal substructure a problem can be broken down into simpler subproblems and for each of these subproblems, a solution can be calculated and the result can be stored in a

table, so that a solution has to be found only once. When a subproblem reoccurs, the cached result can be returned to save computation time.

On the other hand the storage consumption will expand due to the need for storing the results. But in times of rapidly increasing memory sizes and shrinking costs this draw off can often be considered as negligible.

b) Bellman equation:

The Bellman equation, often also referred to as the Dynamic Programming equation, is a necessary condition for optimality. It states a relationship between the value of a decision problem at a certain point in time depending on the outcome of the initial choices and the value of the remaining decisions that result from the previous choices. This breaks the problem down into subproblems that can be calculated separately, as Bellman's principle of optimality suggests.

Before coming to the actual Bellman equation, some underlying concepts have to be defined. At first, optimizing a problem is always under a certain objective: maximizing utility, minimizing path length, maximizing profits. The function describing these objectives is called the **objective function**. Second, Dynamic Programming needs to keep track of the **state** of the system when evolving over time, because when breaking down a multi-period planning problem into simpler steps at different points in time, in every step the information about the current situation that is needed has to be available. Furthermore, the next state is always affected by some factors in addition to the current **control**. That means, that choosing the control variables may be equivalent to choosing the action to reach the next state. This, in fact means making a decision. In dynamic programming, an optimal plan can now be described by finding a rule that makes an interconnection between what the controls should be when given any possible value of the state. To now achieve the best possible value of the objective, the optimal decision rule has to be found.

Finally, by definition, the optimal decision rule is the one that achieves the best possible value of the objective. When written as a function of the state, the best possible value of the objective is called the **value function**.

Richard Bellman showed, that when calculating a relation between the value function of one step and the value function of the next step, a dynamic optimization problem can be solved in a recursive manner. The relationship connecting these two value functions is called the **Bellman equation** and uses an approach called **backwards induction**.

In this approach, the optimal policy of the last time period t_N is specified at first as a function depending on the value of this state and the optimal value of the objective function. Next, the previous step in time t_{N-1} is calculated by optimizing the current (t_{N-1}) objective function and the optimal value for the future (t_N) objective function. This can be applied recursively until the initial state ($t = 0$) is reached and the first period decision rule is derived, as a function of the initial state and the value, by optimizing the sum of the objective function at the initial time step t_0 and the value of the second timestep t_1 ,

which is dependent on all future decisions. Resulting from that, each periods decision is made by explicitly acknowledging that all future decisions will be optimally made.

So, when used on a dynamic decision problem the function describing the previously stated thoughts is

$$V(x_0) = \max_{a_0} \{F(x_0, a_0) + V(x_1)\} \quad (7)$$

whereas x_t is the state of the system at time t , beginning at $t = 0$, x_0 with $a_0 \in \Gamma(x_0)$ being the action depending on the set of possible actions Γ at t_0 . The transition from one state x to another, is described by $T(x, a)$, when action a is chosen at state x , and the payoff of taking a in state x is $F(x, a)$. The Bellman equation can be simplified even more when dropping the time dependency resulting in a general functional equation

$$V(x) = \max_{a \in \Gamma(x)} \{F(x, a) + V(T(x, a))\} \quad (8)$$

Solving this functional equation means now finding the unknown function V , which is the value function of the Problem.

III. FINDING THE VALUE FUNCTION AND IMPLEMENTING THE MAZE

A. Finding the value function

Remember the initial task where you are given a grid map $M = \mathbb{N}^{n \times m}$ containing static walls $M(i, j) = \infty$ representing the maze, and exit(s) of the labyrinth which the player wants to reach.

The goal of this paper is to point out ways to determine an optimal policy $P = \mathbb{N}^{n \times m} \rightarrow A$ to exit the labyrinth by using Dynamic Programming. To achieve this, a way to calculate the value $v(i, j)$ of the current position is searched from what the actual best action can be derived.

To do so, we have to solve the Bellman equation for the task of finding all pair shortest paths in the maze.

This means, that for any arbitrary point in the maze a value function

$$v(i, j) \in \mathbb{N} \quad (9)$$

representing the length of the shortest path of the position (i, j) to the nearest goal in the maze has to be found.

Therefore let $c(i, j) \in \mathbb{R}$ be the cost it would take a prisoner in the maze to step on the field (i, j) .

So now let us at first begin with the trivial cases.

If a prisoner is already standing on a goal field, the cost to reach the goal would, obviously, be

$$c(i, j) = 0, \quad \text{if } (i, j) \in \{X_{\text{goal}}\} \quad (10)$$

. Furthermore, also the value of a goal would be 0, because the length of the shortest path from a position to itself is 0. In addition, if the prisoner isn't even in the maze, he will never be able to reach a goal, resulting in

$$v(i, j) = \begin{cases} \infty, & \text{if } i < 1 \text{ or } i > n \text{ or } j < 1 \text{ or } j > m \\ 0 & \text{if } (i, j) \in \{X_{\text{goal}}\} \end{cases} \quad (11)$$

. At last if the position we're looking for is a wall, the cost is, by definition

$$c(i, j) = \infty, \quad \text{if } (i, j) \in \{\text{Walls}\} \quad (12)$$

and the value of the tile is also ∞ , because there will never be a payoff of stepping into a wall, because it isn't possible at all.

Now to the non trivial cases.

Therefore let us have a look at a maze looking like

1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
1	1	1	1	1

All tiles are having a cost of 1.
 $c(i, j) = 1, \quad \forall (i, j) \in n \times m$

Using the approach of backwards induction introduced by Richard Bellman we can start from the goal state in the grid and calculate the value at the step $N : v(3, 4) = 0$. Now having the result for N , we can induce backwards the results for the step $(N - 1)$, whereas their results are

$$v(2, 4) = v(3, 4) + c(3, 4) = 0 + 1 = 1 \quad (13)$$

$$v(4, 4) = v(3, 4) + c(3, 4) = 0 + 1 = 1 \quad (14)$$

$$v(2, 3) = v(3, 4) + c(3, 4) = 0 + 1 = 1 \quad (15)$$

$$v(2, 5) = v(3, 4) + c(3, 4) = 0 + 1 = 1 \quad (16)$$

When exemplary applying this procedure at the position $v(2, 4)$, we can again give the equations for the positions lying around it

$$v(1, 4) = v(2, 4) + c(2, 4) = 1 + 1 = 2 \quad (17)$$

$$v(3, 4) = v(2, 4) + c(2, 4) = 1 + 1 = 2 \quad (18)$$

$$v(2, 3) = v(2, 4) + c(2, 4) = 1 + 1 = 2 \quad (19)$$

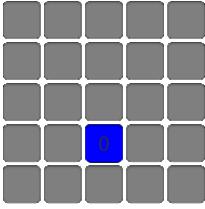
$$v(2, 5) = v(2, 4) + c(2, 4) = 1 + 1 = 2 \quad (20)$$

Now we have calculated the values of the step $(N - 2)$ depending on the decisions that will be made optimally in $(N - 1)$ and N . The results are getting stored in a second matrix $V^{n \times m}$ which holds the results of the value function. So $V(3, 4) = 0, V(2, 4) = 1 \dots$

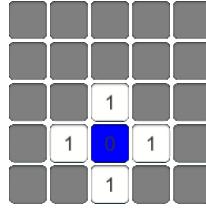
One thing that has to be considered at this step is, that the value for the tile $v(3, 4)$ is getting calculated again, which might cause troubles due to the fact that now the value of this position would be

$$v(3, 4) = v(2, 4) + c(2, 4) = (v(3, 4) + c(3, 4)) + c(2, 4) = 2 \quad (21)$$

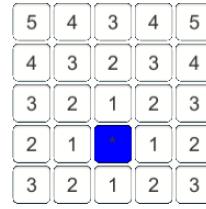
which is in fact bigger than the initially calculated value of 0. At this point we have to take into account if there is already a result of the value function memoized in the matrix V . If there is a result, and the value is less than the one currently calculated, the result can be discarded.



Only the value of the goal position is known at step N .



At step $(N - 1)$ the values of the four adjacent positions are calculated.



At step $(N - 5)$ the values of all the positions are calculated.

Following this procedure, a general relation between the steps can be stated

$$v(i-1, j) = v(i, j) + c(i, j) \quad \text{if } v(i-1, j) < V(i, j) \quad (22)$$

$$v(i+1, j) = v(i, j) + c(i, j) \quad \text{if } v(i+1, j) < V(i, j) \quad (23)$$

$$v(i, j-1) = v(i, j) + c(i, j) \quad \text{if } v(i, j-1) < V(i, j) \quad (24)$$

$$v(i, j+1) = v(i, j) + c(i, j) \quad \text{if } v(i, j+1) < V(i, j) \quad (25)$$

With this equation, the value of every position in the maze can be calculated, which leads to the function solving the functional Bellman equation

$$v(i, j) = \begin{cases} \infty & \text{if } (i, j) \in \{\text{Walls}\} \cup \{\mathbb{R}^2 \setminus \{n \times m\}\} \\ 0 & \text{if } (i, j) \in \{\text{Goals}\} \\ \min\{v(i, j+1) + c(i, j+1), \\ v(i+1, j) + c(i+1, j), \\ v(i, j-1) + c(i, j-1), \\ v(i-1, j) + c(i-1, j)\} & \text{if in the maze and not a wall or a goal} \end{cases} \quad (26)$$

B. Finding optimal policies



Starting from the red tile an optimal policy will be found.



So looking at the adjacent tiles the green one is one of the tiles with the lowest value.



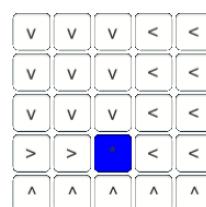
Following equation 27 a path is getting found.



The path is reaching closer to the goal.



The path has reached the goal and the algorithm terminates returning one shortest path from red to blue.



When 27 is applied to every tile, the best action at every position can be given.

With the value of every tile calculated and stored in the matrix V , it is now trivial to find the optimal policy for any

position.

Consider being at any position x in the maze. All you can do is move into one of the four directions right, down, left and up.

Due to the fact that we are searching for the shortest path to the nearest goal, and the value of every tile is known, we can now make our decision where to step depending on the values of the tiles adjacent to the current position. By always choosing the tile with the lowest value, the best action $a \in A$ can be chosen accordingly.

Defined in a mathematical sense, the equation describing the value of every tile is the equation

$$a(i, j) = \begin{cases} \text{up} & \text{if } v(i, j-1) = \min_{(x,y)} \{v(x, y) + c(x, y)\} \\ \text{right} & \text{if } v(i+1, j) = \min_{(x,y)} \{v(x, y) + c(x, y)\} \\ \text{down} & \text{if } v(i, j+1) = \min_{(x,y)} \{v(x, y) + c(x, y)\} \\ \text{left} & \text{if } v(i-1, j) = \min_{(x,y)} \{v(x, y) + c(x, y)\} \end{cases} \quad (27)$$

$$\text{With } (x, y) \in \{(i+1, j), (i-1, j), (i, j+1), (i, j-1)\} \quad (28)$$

IV. RELATED WORK

A. Robot Motion Planning

But how can the Dynamic Programming approach be used for robot motion planning with more than 2 dimensions? Imagine a two-armed robot having six degrees of freedom for every arm resulting in a twelve-dimensional motion planning problem.

For every possible configuration of the arms a node in a graph G can be defined holding the specific attributes of that configuration. For every transition between two configurations an edge between these two can be drawn describing uniquely the transformation between the two states.

With this approach every N -dimensional problem can be projected down onto a two-dimensional plane where the algorithms can easily be applied.

Based on this method many robot motion planning algorithms have been introduced.

B. Algorithmic Design of Feasible Trajectories

Written by Steven M. LaValle, this paper summarizes the recent development of algorithms that construct feasible trajectories made at the University of Illinois. They use et al dynamic programming approaches that produce approximately-optimal solutions for low-dimensional problems.

The difficulties they had to come over are on the one side differential constraints, which means that in a systems the number of actuators is strictly less than the dimension of the configuration space. For example imagine a spacecraft floating in \mathbb{R}^3 that has three thrusters. These are called the actuators of the spacecraft system. Opposite, the spacecraft has six degrees of freedom (X-Pos., Y-Pos., Z-Pos., Yaw, Pitch, Roll) which makes it an underactuated system having differential constraints.

On the other hand they had to deal with global constraints arising from robot collisions.

1) problem formulation:

- 1) State Space: $X \subset \mathbb{R}^n$
- 2) Boundary Values: $x_{\text{init}} \in X, X_{\text{goal}} \subset X$
- 3) Constraint Satisfaction: A function, $D : X \rightarrow \{\text{true, false}\}$, that determines whether global constraints are satisfied from state x .
- 4) Inputs: $U(x) \subset \mathbb{R}^m, \forall x \in X$, which specifies the set of controls.
- 5) State Transition Equation: $\dot{x} = f(x, u)$

2) Dynamic Programming:

Using Bellman's principle of optimality, let k refer to a *stage* or *time step*, and let x_k and u_k refer to the state and input at stage k . For their purpose $K + 1$ represented the final stage. The cost function used to describe this problem is of the form

$$L(x_1, \dots, x_{K+1}, u_1, \dots, u_K) = \sum_{k=1}^K l(x_k, u_k) + l_{K+1}(x_{K+1}) \quad (29)$$

whereas $l_{K+1}(x_{K+1}) = 0$ if $x_{K+1} \in X_{\text{goal}}$, and $l_{K+1}(x_{K+1}) = \infty$ else.

They used classical Cost-to-go Iterations that can solve this problem numerically, whereas the cost-to-go function they used is

$$L_k^*(x_k) = \min_{u_k, \dots, u_K} \left\{ \sum_{i=k}^K l(x_i, u_i) + l_{K+1}(x_{K+1}) \right\}. \quad (30)$$

This now induces the calculation of the cost-to-go functions iteratively from stage $K + 1$ to 1. In each iteration, L_k^* is computed using the result of L_{K+1}^* by using the following dynamic programming equation, which involves a local optimization over the inputs:

$$L_k^*(x_k) = \min_{u_k} \{l(x_k, u_k) + L_{K+1}^*(x_{K+1})\}. \quad (31)$$

For further reading see [8]

C. Rapidly exploring Random Trees (RRT)

1) What is a RRT:

Another very interesting field of research are Rapidly-exploring Random Trees, which were first introduced by LaValle in October 1998 [?].

The main idea behind RRTs is a randomized data structure that is designed for a broad class of path planning problems that are viewed as search in a metric space X for a continuous path from an initial state x_{init} to a goal region $X_{\text{goal}} \subset X$

WHAT IS A METRIC

In addition, a fixed obstacle region $X_{\text{obs}} \subset X$ must be avoided, but there is no explicit representation of it available, thus one can only check if a specific state lies in X_{obs} .

WHAT IS EXPLICIT REPRESENTATION

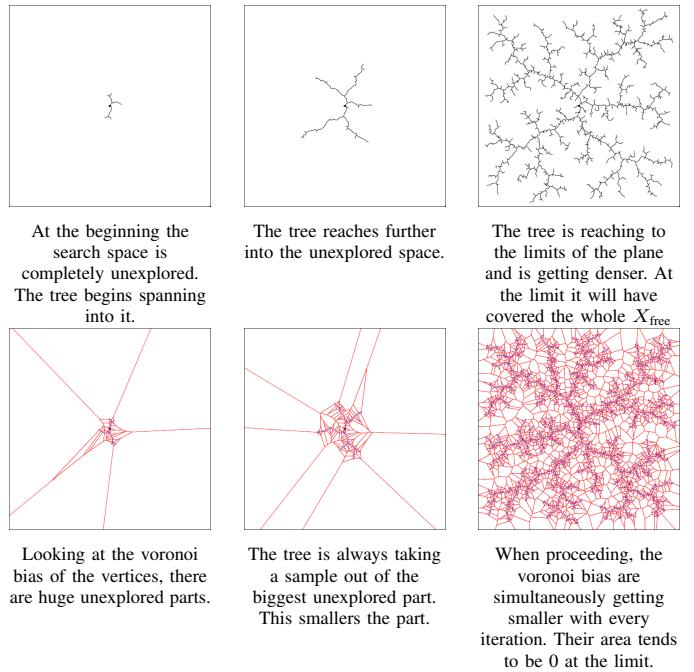
The nice thing about RRTs is, that X_{obs} not only represents solid obstacles like walls, trees or blocks as known from computer games, but can represent any type of constraints. These include velocity bounds of a self-driving car, configurations at which a robot is in collision with an obstacle in the world, Planning Detailed Animation at

runtime and any other constraints occurring in motion planning problems.

When constructing the RRT, the vertices represent the states of the system in X_{free} , and each edge will correspond to a path that lies in completely in X_{free} . [?]

2) How does a RRT look like and how to construct:

- 1) Pick a random sample in the search space.
- 2) Find the nearest neighbor of that sample.
- 3) Select an action from the neighbor that heads towards the random sample.
- 4) Create a new sample based on the outcome of the action applied to the neighbor.
- 5) Add the new sample to the tree, and connect it to the neighbor.



3) RRTs and dynamic programming:

In the dissertation of Oktay Arslan at the Georgia Institute of Technology in December 2015, he gives an extension of the solely sampling based approach of the original RRT.

He introduces an algorithm called RRT[#] which utilizes the main ideas from RRT, but expand it to a relaxation step, used in Dynamic Programming.

Relaxation is a term used iterative methods for solving systems of equations, in this case of optimization problems in numerical mathematics.

They can be written down in the form of

$$x(k+1) = f(x(k)), \quad k = 0, 1, \dots \quad (32)$$

where $x(k) \in \mathbb{R}^n$ is an n-dimensional vector, and $f : \mathbb{R}^n \mapsto \mathbb{R}^n$ is a vector-valued function.

The Bellman function is of this form, what spans the bow between Dynamic Programming and RRT[#].

For understanding the algorithm, some concepts have to be defined:

- Successor vertices: Given a vertex $v \in V$ in a directed Graph $G = (V, E)$, the function $\text{succ} : (G, v) \mapsto V' \subseteq V$ returns the vertices in V that can be reached from vertex v .

$$\text{succ}(G, v) := \{u \in V : (v, u) \in E\}$$

- Predecessor vertices: Given a vertex $v \in V$ in a directed Graph $G = (V, E)$, the function $\text{pred} : (G, v) \mapsto V' \subseteq V$ returns the vertices in V from which vertex v can be reached.

$$\text{pred}(G, v) := \{u \in V : (u, v) \in E\}$$

Following thereof, they rewrote the bellman equation for their purpose in terms of the cost-to-come value as follows:

$$g^*(v_i) = \begin{cases} 0 & v_i = x_{\text{init}} \\ \min_{v_j \in \text{pred}(G, v_i)} (g^*(v_j) + c(v_j, v_i)) & v_i \in V \setminus \{x_{\text{init}}\} \end{cases} \quad (33)$$

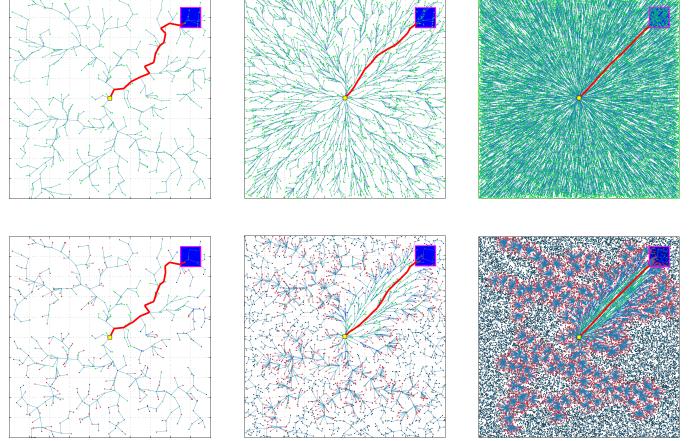
P.42 Based on their formulation of the Bellman equation, they split the Algorithm into two parts:

a) *Exploration*: This samples randomly a point in X_{free} and then extends the underlying graph towards the sampled point by including it as a new vertex in the current graph by connecting the missing edges. [10]

b) *Exploitation*: This step implements the task of improving the cost-to-come values of the current vertices as new information becomes available. It also encodes the lowest-cost path information for the promising vertices (see equation (20) below) and v_{goal}^* (the goal vertex that has the lowest cost-to-come value in the goal set) as a spanning tree rooted at the initial vertex. Cost-to-come values of nonstationary vertices are updated in an order based on their f-values, i.e., an underestimate of the cost of the optimal path from the initial vertex to the goal set passing through the vertex of interest, and ties are broken in favor of vertices which have smaller g-values at each iteration of the Gauss-Seidel version of the Bellman-Ford algorithm.

The following images show the two algorithms.

As you can see, dynamic programming is a great way to improve the runtime and the speed in which an algorithm converges to an optimal solution. The interested reader can have a look at [10] for further information and mathematical proofs.



V. CONCLUSION

As we have seen in this paper, Dynamic Programming is a very straight forward approach for solving multistage decision problems. By breaking down difficult problems into sets of overlapping subproblems with optimal substructure an overall optimal solution can be found by only solving the small subproblems. In robot motion, dynamic programming can be applied in a wide field of tasks and can be a really useful approach that takes the place of the classic algorithms. Many algorithms are already getting improved or replaced by Dynamic Programming algorithms and the set of fields where it can be applied is nearly unlimited.

REFERENCES

- [1] R. Bellman, “On the theory of dynamic programming,” The RAND Corporation, 1952.
- [2] ——, “On a routing problem,” vol. 16, pp. 87–90, 1958. [Online]. Available: <http://www.jstor.org/stable/43634538>
- [3] ——, “The theory of dynamic programming,” The RAND Corporation, 1954. [Online]. Available: <https://www.rand.org/content/dam/rand/pubs/papers/2008/P550.pdf>
- [4] S. P. Bradley, A. C. Hax, and T. L. Magnanti, *Applied mathematical programming*, [19. dr.] ed. Addison-Wesley, 1992, hax, Arnaldo C. (VerfasserIn) Magnanti, Thomas L. (VerfasserIn).
- [5] K. L. Cooke and E. Halsey, “The shortest route through a network with time-dependent internodal transit times,” vol. 14, pp. 493–498, 1966. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/0022247X66900096>
- [6] “Operations research proceedings 1994: Selected papers of the international conference on operations research, berlin, august 30 – september 2, 1994,” 1995.
- [7] S. Dreyfus, “Richard bellman on the birth of dynamic programming,” vol. 50, pp. 48–51, 2002. [Online]. Available: <https://doi.org/10.1287/opre.50.1.48.17791>
- [8] S. LaValle, “From dynamic programming to rrt: Algorithmic design of feasible trajectories,” University of Illinois. [Online]. Available: <http://msl.cs.illinois.edu/lavalle/papers/Lav02.pdf>
- [9] S. M. LaValle, *Planning Algorithms*. Cambridge University Press, 2006 % References from the book, available at <http://planning.cs.uiuc.edu/>.
- [10] O. Arslan, “Machine learning and dynamic programming algorithms for motion planning and control,” 2015-12-01. [Online]. Available: <http://hdl.handle.net/1853/54317>
- [11] P. D. Pra, C. Rudari, and W. J. Runggaldier, “On dynamic programming for multistage decision problems under uncertainty,” in *Operations Research Proceedings 1994: Selected Papers of the International Conference on Operations Research, Berlin, August 30 – September 2, 1994*, U. Derigs, A. Bachem, and A. Drexl, Eds. Springer Berlin Heidelberg, 1995, pp. 70–75. [Online]. Available: