

Introduction to Deep Reinforcement Learning

Jill-Jênn Vie

November 14, 2025

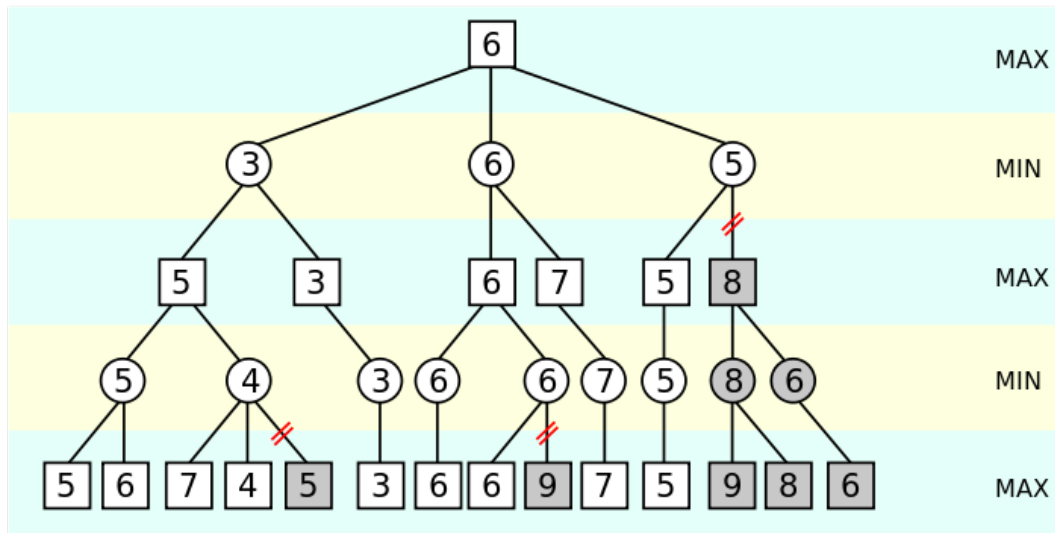
Reinforcement learning

- ▶ Train an agent from interaction with an environment
- ▶ Actions from the agent impact the data collected

Applications

- ▶ Robotics, self-driving cars
- ▶ Games
- ▶ Recommender systems
- ▶ LLMs like ChatGPT

Games on trees or DAGs: minimax (von Neumann, 1928; McCarthy, 1958)



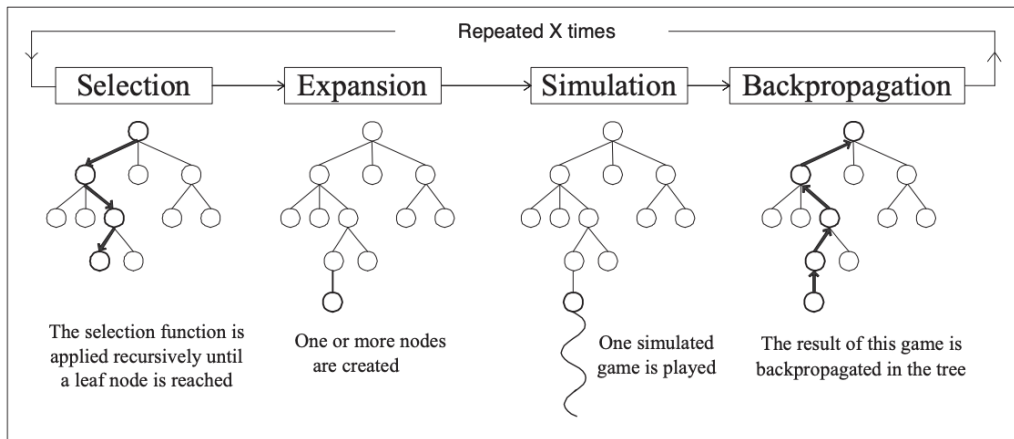
Monte Carlo Tree Search (coined by Rémi Coulom, Inria Lille, in 2006)

1987: Bruce Abramson, minimax search with an expected-outcome model

2012: Cameron Browne et al. "A survey of monte carlo tree search methods."

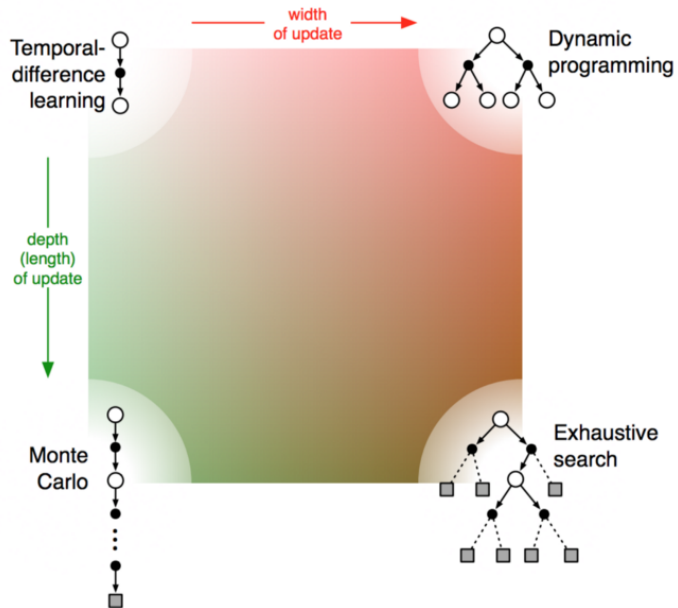
2015: AlphaGo combines MCTS and deep learning and beats professional human Go player

Image source: Chaslot et al. (2008)



Algorithm/Discovery	Nature Paper Title	Core Breakthrough
Atari (DQN)	<i>Human-level control through deep reinforcement learning</i> (2015)	First time an agent combined deep neural networks with Q-learning (a model-free method) to achieve human-level performance across a wide range of visually complex tasks (Atari games).
AlphaGo	<i>Mastering the game of Go with deep neural networks and tree search</i> (2016)	First AI to defeat a top professional human Go player. It combined Supervised Learning (pre-training on human games) with Reinforcement Learning (self-play) and Monte Carlo Tree Search (MCTS) .
AlphaGo Zero	<i>Mastering the game of Go without human knowledge</i> (2017)	Eliminated human data. Trained purely through self-play (tabula rasa), achieving stronger performance than AlphaGo. Simplified the network architecture and used a single neural network for both policy and value.
MuZero	<i>Mastering Atari, Go, Chess and Shogi by planning with a learned model</i> (2020)	Eliminated knowledge of environment rules. Masters board games and Atari games by learning a compressed model of the environment's dynamics, focusing only on predicting reward, value, and policy —the elements essential for planning.

Various aspects of RL (Sutton & Barto)



In RL, episodes are sampled according to a policy

Agent observes state S_t chooses action A_t

You can control the agent: **policy** $\pi(a|s) = \pi(A_t = a|S_t = s)$

$$S_t \rightarrow A_t \rightarrow R_t \rightarrow S_{t+1} \rightarrow A_{t+1} \rightarrow \cdots \rightarrow R_T$$

This is called one *episode* of learning.

In RL, episodes are sampled according to a policy

Agent observes state S_t chooses action A_t

You can control the agent: **policy** $\pi(a|s) = \pi(A_t = a|S_t = s)$

Environment replies with reward R_t and new state S_{t+1}

You cannot control the environment: $p(s', r|s, a) = p(S_{t+1} = s', R_t = r|S_t = s, A_t = a)$

$$S_t \rightarrow A_t \rightarrow R_t \rightarrow S_{t+1} \rightarrow A_{t+1} \rightarrow \cdots \rightarrow R_T$$

This is called one *episode* of learning.

In RL, episodes are sampled according to a policy

Agent observes state S_{t+1} chooses action A_{t+1}

You can control the agent: **policy** $\pi(a|s) = \pi(A_t = a|S_t = s)$

Environment replies with reward R_t and new state S_{t+1}

You cannot control the environment: $p(s', r|s, a) = p(S_{t+1} = s', R_t = r|S_t = s, A_t = a)$

$$S_t \rightarrow A_t \rightarrow R_t \rightarrow S_{t+1} \rightarrow \textcolor{red}{A_{t+1}} \rightarrow \cdots \rightarrow R_T$$

This is called one *episode* of learning.

Objectives

Let $\gamma \in (0, 1)$ be a discount factor

The goal is to optimize the *discounted return*:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad \text{where } T \text{ can be } \infty$$

The *value function* v_π for a policy π is $v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$

The *action-value function* q_π for a policy π is $q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$

What can be done with these functions?

Policy evaluation

Given policy π , compute value v_π

$$v_\pi(s) = \mathbb{E}_a q_\pi(s, a) = \sum_a \pi(a|s) q_\pi(s, a)$$

Policy improvement

Given value v_π on states, improve policy π

i.e. pick best action a : $\operatorname{argmax}_a q_\pi(s, a)$

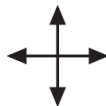
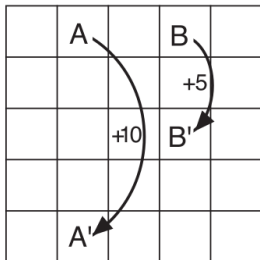
$$q_\pi(s, a) = \mathbb{E}_\pi[r + \gamma v_\pi(s')] = \sum_{r, s'} p(r, s'|s, a)[r + \gamma v_\pi(s')]$$

A first example: policy evaluation on Gridworld

$$S \in \{(i, j)\}_{1 \leq i, j \leq 5}$$

$$A \in \{\leftarrow, \rightarrow, \uparrow, \downarrow\}$$

$$R = \begin{cases} 10 & \text{if leaving from A} \\ 5 & \text{if leaving from B} \\ -1 & \text{if hitting a wall} \\ 0 & \text{otherwise.} \end{cases}$$



Actions

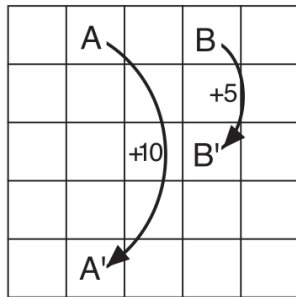
3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

Gridworld

$$\pi_{\text{uniform}} = (1/4, 1/4, 1/4, 1/4)$$

$v_{\pi_{\text{uniform}}}$

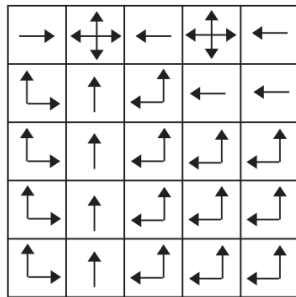
A first example: policy improvement on Gridworld



Gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

v_*



π_*

Richard Bellman's principle of optimality (1952) for dynamic programming

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

Bellman's equation

$$\underbrace{V(s)}_{\text{value}} = \max_a \underbrace{R(s, a)}_{\text{reward}} + \gamma \underbrace{V(T(s, a))}_{\text{transition}}$$

$$v_{\pi^*}(s) = \max_a \mathbb{E}_{\pi^*}[r + \gamma v_{\pi^*}(s')]$$

Rings a bell?

Richard Bellman's principle of optimality (1952) for dynamic programming

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

Bellman's equation

$$\underbrace{V(s)}_{\text{value}} = \max_a \underbrace{R(s, a)}_{\text{reward}} + \gamma \underbrace{V(T(s, a))}_{\text{transition}}$$

$$v_{\pi^*}(s) = \max_a \mathbb{E}_{\pi^*}[r + \gamma v_{\pi^*}(s')]$$

Rings a bell?

Bellman-Ford in shortest paths

$$V(u) = \max_v -w_{u,v} + V(v) \quad V = -d \quad \gamma = 1$$

Bellman-Ford

```
const int oo = 1e9;
int n, m;
int dist[N];
vector<tuple<int, int, int>> edge;
```

```
bool bellmanford (int u0) {
    fill_n (dist, n, +oo);
    dist[u0] = 0;
    bool stable = false;
    for (int t = 0; t < n && !stable; t++) {
        stable = true;
        for (auto[u, v, c] : edge) if (dist[u] < +oo && dist[u] + c < dist[v]) {
            dist[v] = dist[u] + c;
            stable = false;
        }
    }
    return stable;
}
```

Repeat at most $|V|$ times

For each edge $(u, v) \in E$

$$d(v) = \min(d(v), d(u) + w_{uv})$$

$d_k[v]$ = shortest length from source to v using at most k edges

Value iteration

Value Iteration, for estimating $\pi \approx \pi_*$

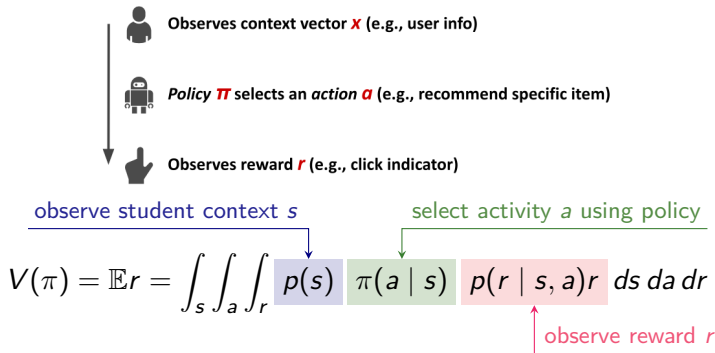
Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

```
|  $\Delta \leftarrow 0$   
| Loop for each  $s \in \mathcal{S}$ :  
|    $v \leftarrow V(s)$   
|    $V(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$   
|    $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
until  $\Delta < \theta$ 
```

Output a deterministic policy, $\pi \approx \pi_*$, such that
$$\pi(s) = \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

Particular case: $\gamma = 0$ is the setting of contextual bandits



Online learning of classifiers: if the classifier looks at an image x picks a class y it may receive a reward of 1 if it picked the right class, -1 otherwise

ChatGPT: given a prompt x it should choose an answer y and gets a reward $r(x, y)$ that it estimates itself from preferences

Reinforcement learning from human feedback: InstructGPT, ChatGPT

1. Collect demonstration data, and train a supervised policy $\pi_0(y|x)$ (based on GPT-3)
2. Collect comparison data, train a reward model (Elo rating from “only” 50k pairs)

$$\text{loss}(\theta) = -\mathbb{E}_{(x, y_w, y_\ell) \sim D} \log \underbrace{\frac{\sigma(r_\theta(x, y_k) - r_\theta(x, y_\ell))}{\Pr(\text{"answer } y_k \text{ is preferred to } y_\ell")}}$$

3. Optimize a policy against the reward model using PPO.

$$\text{objective}(\phi) = \mathbb{E}_{(x, y) \sim \pi_\phi} r_\theta(x, y) - \beta \text{KL}(\pi_\phi, \pi_0)$$

Ouyang, Wu, Jiang, Almeida, Wainwright, Mishkin, ... & Lowe (NeurIPS 2022). Training language models to follow instructions with human feedback.

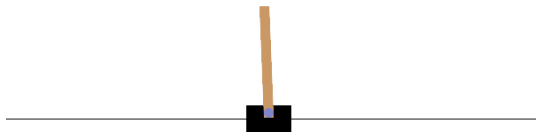
A recent paper suggests to drop part 3 (RL) and focus on part 2 (supervised reward model):

Rafailov, Sharma, Mitchell, Ermon, Manning and Finn (arXiv 2023.06).

Direct Preference Optimization: Your Language Model is Secretly a Reward Model.

Practical

CartPole-v1

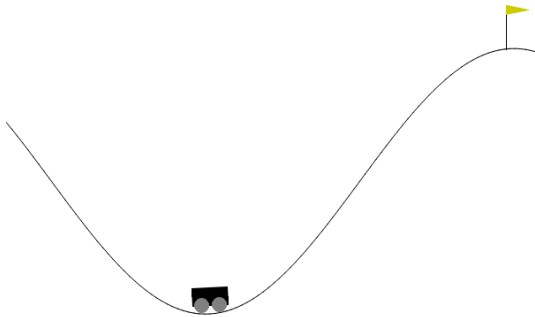


$$S = (x, \dot{x}, \theta, \dot{\theta}) \in \mathbb{R}^4$$

$$A \in \{\leftarrow, \rightarrow\}$$

$$R = 1$$

MountainCar-v2

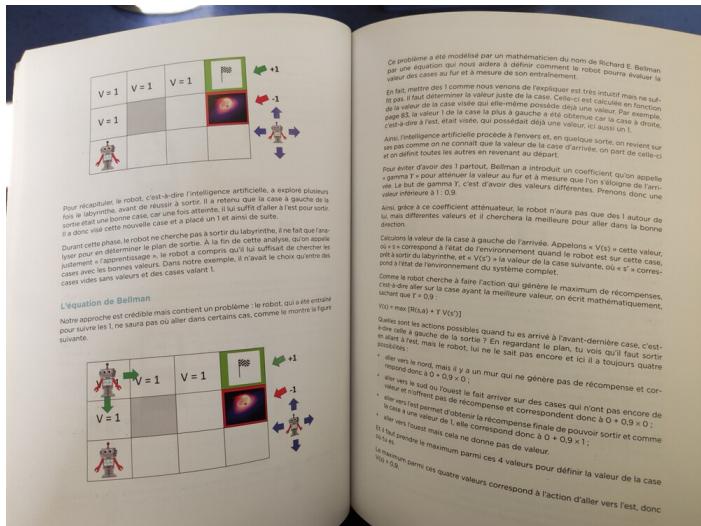


$$S = (x, \dot{x}) \in \mathbb{R}^2$$

$$A \in \{\leftarrow, 0, \rightarrow\}$$

$$R = -1$$

References

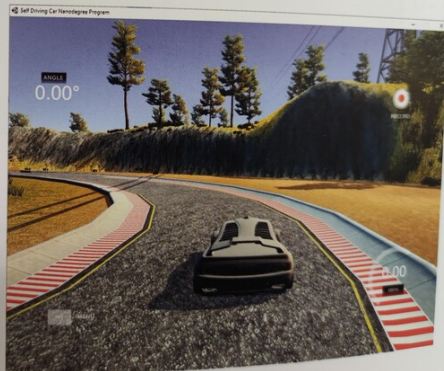


Code une voiture autonome avec Python

Dans cet exercice, tu vas jouer en ligne et entraîner une voiture autonome. Comment vas-tu procéder ? Avant le début de la course sur le circuit, tu pourras activer un enregistrement qui récupérera les multiples images de la vidéo de ton jeu qui serviront à générer un modèle basé sur un réseau de neurones à convolution. Te rappelles-tu des réseaux spécialisés pour la reconnaissance d'images ?

Nous nous appuyerons sur les outils développés par NVIDIA, une société américaine spécialisée dans la fabrication de cartes graphiques qui développe depuis 2017 de nouvelles cartes pour l'intelligence artificielle.

Les voitures autonomes sont l'un des domaines de recherche et d'activité les plus dynamiques. Ce qui semblait relever de la



This basic equation can be used to study optimal trajectories, including launching, geodesics, paths of minimum length through discrete sets, the r -th shortest paths through networks, soft landing on the moon and other planets, chemical process control, reactor control, and many other problems and problems.

Questions of computational solution will be discussed in the next chapter.

4.17 The Brachistochrone

As an example of a problem in the calculus of variations which can be explicitly resolved, either by classical techniques, or by means of the functional equation technique which we shall pursue here, let us consider the problem of determining a curve connecting the two points P and Q with the property that a particle sliding along this curve under the sole influence of gravity will arrive at Q in a minimum time.

Without loss of generality, let P be the origin, and let the axes be oriented as indicated below.

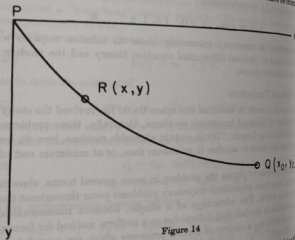


Figure 14

Omitting the gravitational constant, the problem is that of minimizing the functional

$$(4.56) \quad J(y) = \int_0^{x_0} \left(\frac{1 + y'^2}{y} \right)^{1/2} dx,$$

over all curves satisfying the end conditions

$$(4.57) \quad y(0) = 0, \quad y(x_0) = y_0$$

To treat the problem by dynamic programming techniques, we let $f(x, y)$ denote the minimum time required to go from the generic point $R(x, y)$ to (x_0, y_0) . The equation in (4.56) is in this case

$$(4.58) \quad f(x, y) = \min_v \left[\left(\frac{1 + y'^2}{y} \right)^{1/2} \Delta + f(x + \Delta, y + y' \Delta) + O(\Delta^2) \right]^*.$$

Passing to the limit as $\Delta \rightarrow 0$, this yields the nonlinear partial differential equation

$$(4.59) \quad 0 = \min_v \left[\left(\frac{1 + y'^2}{y} \right)^{1/2} + f_x + y' f_y \right],$$

equivalent to the two equations

$$(4.60) \quad \begin{aligned} (a) \quad 0 &= \left(\frac{1 + y'^2}{y} \right)^{1/2} + f_x + y' f_y, \\ (b) \quad 0 &= \frac{y'}{[y(1 + y'^2)]^{1/2}} + f_y. \end{aligned}$$

To eliminate the function f , and thus obtain an equation for y' , we differentiate (4.60a) with respect to x and take the partial derivative of (4.60a) with respect to y , obtaining

$$(4.61) \quad \begin{aligned} (a) \quad \frac{d}{dx} \left[\frac{y'}{[y(1 + y'^2)]^{1/2}} \right] + f_{xx} + f_{xy} y' &= 0, \\ (b) \quad -\frac{1}{2} (1 + y'^2)^{-3/2} y' + f_{yy} y' &= 0. \end{aligned}$$

Thus

$$(4.62) \quad \frac{d}{dx} \left[\frac{y'}{[y(1 + y'^2)]^{1/2}} \right] + \frac{1}{2} (1 + y'^2)^{-3/2} = 0.$$

A first integral is

$$(4.63) \quad \left(\frac{1 + y'^2}{y} \right)^{1/2} - \frac{y'^2}{[(1 + y'^2)y]^{1/2}} = k,$$

or

$$(4.64) \quad \frac{1}{[y(1 + y'^2)]^{1/2}} = k.$$

This is equivalent to Snell's law for the propagation of light. It is interesting to interpret this well-known physical principle as an optimal policy in a multistage decision process.

The integration can now be readily completed to obtain the standard representation of a brachistochrone in parametric form:

$$(4.65) \quad \begin{aligned} y &= (1 - \cos t)/2k^2, \\ x &= c_1 + (t - \sin t)/2k^2. \end{aligned}$$

* We have deliberately kept y' as the slope, instead of v , in order to manipulate the Euler equation.

The great drawback of dynamic programming is, as Bellman himself calls it, the "curse of dimensionality." Even recording the solution to a moderately complicated problem involves an enormous amount of storage. If we want only one optimal path from a known initial point, it is wasteful and tedious to find a whole field of extremals; if we need feedback, a perturbation feedback scheme is often quite adequate (see Chapter 6, Neighboring Extremals and the Second Variation).

Derivation of the Euler-Lagrange equations from the Hamilton-Jacobi equation. Consider a particular optimal path and its associated optimal control function. Then we have

$$\frac{d\lambda^T}{dt} = \frac{d}{dt} \left(\frac{\partial J^0}{\partial \dot{x}} \right) = \frac{\partial^2 J^0}{\partial x^2} \dot{x} + \frac{\partial^2 J^0}{\partial x \partial t} \quad (4.2.19)$$

Partial differentiation of (4.2.15) with respect to x , considering $u^0 = u^0(x, t)$, gives

$$\frac{\partial^2 J^0}{\partial x \partial t} + \frac{\partial L}{\partial u} \frac{\partial u^0}{\partial x} + f^T \frac{\partial^2 J^0}{\partial x^2} + \frac{\partial J^0}{\partial x} \left(\frac{\partial f}{\partial x} + \frac{\partial f}{\partial u} \frac{\partial u^0}{\partial x} \right) = 0 \quad (4.2.20)$$

Now the coefficient of $\partial u^0 / \partial x$ in (4.2.20) vanishes on an optimal path according to (4.2.17).† Using (4.2.20) in (4.2.19), then, we obtain

$$\frac{d\lambda^T}{dt} = -\lambda^T \frac{\partial f}{\partial x} - \frac{\partial L}{\partial x} \quad (4.2.21)$$

which, along with (4.2.17), are the Euler-Lagrange equations.

Furthermore, the fact that J^0 is equal to ϕ on $\psi = 0$ implies that there exists a vector ν such that

$$\left. \frac{\partial J^0}{\partial x} \right|_{t_f} = \left(\frac{\partial \phi}{\partial x} + \nu^T \frac{\partial \psi}{\partial x} \right)_{t_f} \triangleq \lambda^T(t_f) \quad (4.2.22)$$

In words, the change in cost due to admissible change in state (that is, $d\psi = 0$) is given by a linear combination of the gradient of ϕ with respect to state and the gradients of ψ (constraints) with respect to state (see Section 1.2).

Combinatorial problems. Dynamic programming is especially useful in solving multistage optimization problems in which there are only a small number of possible choices of the control at each

†When there are inequality constraints on the control variable, it can be shown (e.g., by defining a modified Hamiltonian as in Chapter 3 or see Kalman (1963) or Dreyfus (1965)) that the term $(L_u + f_u^T \lambda) u_u^*$ still vanishes.

stage and in which no derivative information is available. Consider a simple example with only two possible choices of control at each stage. We would like to find the path from A to B in Figure 4.2.1,

