

Business Analytics – 2020/21

Modeling for dynamic programming

Reinforcement learning = model free dynamic programming

\ Algoritmo puramente statistico che NON conosce $T_{i,j}(\cdot)$

Prof. Paolo Brandimarte

Dip. di Scienze Matematiche – Politecnico di Torino

e-mail: paolo.brandimarte@polito.it

URL: <https://staff.polito.it/paolo.brandimarte>

This version: July 27, 2022

NOTE: For internal teaching use within the Masters' Program in Mathematical Engineering. Do not post or distribute.

References

These slides are taken from my book:

P. Brandimarte. *From Shortest Paths to Reinforcement Learning: A MATLAB-Based Introduction to Dynamic Programming*. Springer, 2021.

Other references:

- P. Brandimarte. *An Introduction to Financial Markets: A Quantitative Approach*. Wiley, 2018.
- J.Y. Campbell, L.M. Viceira. *Strategic Asset Allocation*. Oxford University Press, 2002.
- W.B. Powell. *Approximate Dynamic Programming: Solving the Curses of Dimensionality* (2nd ed.). Wiley, 2011.
- G.J. van Ryzin, K.T. Talluri. *An Introduction to Revenue Management*. In INFORMS TutORials in Operations Research, 2005, <https://pubsonline.informs.org/doi/10.1287/educ.1053.0019>.

Modelling for DP

The DP principle is a flexible concept for the decomposition of a multistage decision problem into a sequence of single-stage problems, by balancing the immediate contribution and the expected contributions from future decisions.

This hinges on the value function defined by the DP recursion

Alcune situazioni prevedono lo scambio del valore atteso con l'ottimizzazione
opt 2. Il problema e' così
opt 2. Lo modifichi tu in questo modo

$$V_t(s_t) = \underset{x_t \in \mathcal{X}(s_t)}{\text{opt}} \left\{ f_t(s_t, x_t) + \gamma \mathbb{E}[V_{t+1}(s_{t+1}) | s_t, x_t] \right\}. \quad \begin{array}{l} \text{Equazione ricorsiva della programmazione dinamica a orizzonte finito tipo stocastico} \\ \text{Principio della programmazione dinamica} \end{array} \quad (1)$$

Il contributo immediato puo' essere stocastico \Rightarrow andrebbe dentro al valore atteso

Equation (1) is only one possible form of DP recursion.

→ Notiamo che conosciamo s_t , e non i successivi in generale

We may adopt more specific forms, e.g., when risk factors are discrete random variables and the expectation boils down to a sum.

A more radical rephrasing is sometimes adopted, where we swap expectation and optimization.

This may occur because of the information structure of the problem at hand, or it may be the result of some manipulation aimed at avoiding the solution of a difficult stochastic optimization subproblem. → Si ottengono introducendo gli stati post decisione = post decision state

One well-known case occurs in Q -learning, a form of reinforcement learning where the state value function $V(s)$ is replaced by Q -factors $Q(s, x)$ representing the value of state-action pairs.

↓
Scegliere come rappresentare lo stato
non è ovvio e può facilitare la soluzione

Value function dipendenti dalla coppia
stato - azione . E' complicato, ma c'e' un
grosso vantaggio: non devi info sulle distribuzioni
di probabilità avere

More generally, we will see how the swap of optimization and expectation can be obtained by introducing the concept of post-decision states.

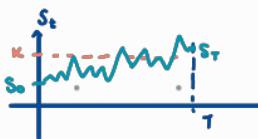
Also finding a suitable description of the system state may require careful modeling. For instance, a redefinition of the state space may be needed to eliminate path-dependencies and to obtain an augmented and equivalent Markovian model, in order to make the problem amenable to a DP approach.

Es. opzione
call asiatica

We should also note that, although it is natural to think of states and control decisions as scalars or vectors, they may consist of different objects, like sets.

The bottom line is that modeling skills come hand in hand with algorithmic knowledge, when tackling a decision problem by DP. These skills can only be honed by experiencing with a diverse array of problems.

Consideriamo l'opzione call \Rightarrow ti dà il diritto, non l'obbligo, di comprare l'azione S_t all'istante $t=T$ al prezzo strike K .



cioè il guadagno quando usi l'opzione call

Il pay-off è $\max(S_T - K, 0)$, che ha quindi un grafico del tipo



da cui il pay-off non è mai negativo

- Un'opzione di stile europeo la puoi esercitare solo alla maturità, per cui apparentemente non c'è nulla da ottimizzare

- Un'opzione di stile americana. la puoi esercitare in un qualunque momento prima della maturità

Quando bisogna esercitare l'opzione in questo caso? Si sfrutta la programmazione dinamica \Rightarrow con pay-off $\max(S_t - K, 0)$

Quanto manca alla maturità dell'opzione inciderà sulla scelta

A volte uno vuole eliminare le dipendenze dal percorso. \rightarrow quello che conta è il prezzo adesso, non come ci sei arrivato soprattutto se si ha un processo Markoviano

La variabile di stato è il prezzo dell'azione in questo momento

- Un'opzione è di stile asiatico se il pay-off è $\max\left(\frac{1}{T} \sum_{t=1}^T S_t, 0\right)$ quindi dipende dal prezzo medio osservato fino a quel momento \Rightarrow il processo non è più Markoviano perché dipende dal percorso

Media dei prezzi a tempo discreto osservati precedentemente

$$\frac{1}{T} \sum_{t=1}^T S_t$$

fino a quel momento \Rightarrow il processo non è più Markoviano perché dipende dal percorso

Per togliere la parte «pendency», cioè non Markoviana, si aggiunge una nuova variabile di stato aumentata tendenzialmente un integrale della variabile di stato

Finite Markov decision processes

→ Consideriamo una catena di Markov

The term **Markov decision process (MDP)** is reserved to problems featuring discrete state and action spaces. → con questi spazi discreti posso applicare questa teoria di base (non applicabile)

The term **action** is often adopted to refer to control decisions in this context. Since states and actions are discrete, they can be enumerated, i.e., associated with integer numbers.

Here, we only consider *finite* MDPs. Hence, even though states may correspond to vectors of a multidimensional space, we will refer to states by a notation like

Problema: tirare fuori le probabilità di transizione è difficile $i, j \in \mathcal{S} = \{1, \dots, N\}$,

where N is the size of the state space. Spazio degli stati che sono contabili

By the same token, we may use a notation like a or a_t to refer to actions, as well as $\mathcal{A}_t(i)$ or $\mathcal{A}(i)$ to denote the set of feasible actions at state $i \in \mathcal{S}$ at time t . We will also denote the whole set of possible actions by \mathcal{A} .

The underlying system is essentially a discrete-time and finite **Markov chain**, where it is customary to represent dynamics by a matrix collecting transition probabilities between pairs of states.

In MDPs we use transition probabilities too, but, since we deal with a (partially) controlled Markov chain, the transition probabilities depend on the selected action:

Probabilità difficilmente calcolabili $\Rightarrow \pi_{t+1}(i, a, j)$ — La probabilità dipende dall'azione che scelgo
— probabilità di transizione per catene di Markov a tempo discreto e spazio degli stati finiti o al più numerico.
is the probability of a transition from state i to state j during time interval $t + 1$, after choosing action a at time t .

Example: stochastic inventory control

Let us consider again the toy stochastic inventory control problem, where demand may assume values within the set $\{0, 1, 2\}$, with probabilities 0.1, 0.7, and 0.2, respectively.

Since there is a constraint on the maximum inventory level, $I_t \leq 2$, the state space is

$$\mathcal{S} = \{0, 1, 2\},$$

since inventory may take values 0, 1, and 2. By the same token, the feasible action sets are

$$\mathcal{A}(0) = \{0, 1, 2\}, \mathcal{A}(1) = \{0, 1\}, \mathcal{A}(2) = \{0\}. \rightarrow \text{Puoi ordinare in base al livello iniziale di magazzino}$$

The transition probability matrices $\Pi(a)$, where each element $\pi_{ij}(a)$ contains the transition probability $\pi(i, a, j)$, are as follows:

Calcoliamo le probabilità di transizione che dipendo dall'azione

$$\Pi(0) = \begin{bmatrix} 1 & 0 & 0 \\ 0.9 & 0.1 & 0 \\ 0.2 & 0.7 & 0.1 \end{bmatrix} \xrightarrow{\substack{\text{Se con 1 petro in magazzino mette} \\ \text{chiedono 1 o 2 (oltre 2)}}} \Pi(1) = \begin{bmatrix} 0.9 & 0.1 & 0 \\ 0.2 & 0.7 & 0.1 \\ \cdot & \cdot & \cdot \end{bmatrix}, \quad \Pi(2) = \begin{bmatrix} 0.2 & 0.7 & 0.1 \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \xrightarrow{\substack{\text{Sulle righe stato corrente} \\ \text{Sulle colonne stato futuro}}} \xrightarrow{\substack{\text{Il magazzino non} \\ \text{può essere già pieno}}} \xrightarrow{\substack{\text{Ordino 2 solo} \\ \text{se il magazzino} \\ \text{era completamente} \\ \text{vuoto}}}$$

Here, the dot \cdot is used for rows corresponding to states in which an action is not feasible. For instance, the action $a = 2$ is only feasible at state $i = 0$, because of the constraint on maximum inventory.

La problematica di questo approccio e' quindi che Π non e' facilmente trovabile, soprattutto in casi reali in cui bisogna tenere anche in conto le probabilità di sostituzione

DP recursions

Problemi a tempo finito

In the MDP setting, the value function at any time instant t boils down to a finite-dimensional vector with components $V_t(i)$, and its satisfies

Calcolo delle Value function all'indietro
come sempre. Problemi → spazio degli stati a grandi dimensioni
(...)

$$V_t(i) = \underset{a \in \mathcal{A}(i)}{\text{opt}} \left\{ f_t(i, a) + \gamma \sum_{j \in \mathcal{S}} \pi_{t+1}(i, a, j) V_{t+1}(j) \right\}, \quad i \in \mathcal{S}. \quad (2)$$

↓ azione
 ↓ Stato di partenza
 ↓ programmazione dinamico scontata
 ↓ politica ottimale

In the discounted infinite-horizon case, we have

Ora il problema e' che anche avendo $\pi(j)$, ho bisogno di sapere $f(i, a)$

La Value function dipendono dallo stato

$$V(i) = \underset{a \in \mathcal{A}(i)}{\text{opt}} \left\{ f(i, a) + \gamma \sum_{j \in \mathcal{S}} \pi(i, a, j) V(j) \right\}, \quad i \in \mathcal{S}. \quad (3)$$

↓ La Value function e' un punto fisso dell'operatore
 ↓ opt ⇒ la teoria aiuta a dimostrare ∃!
 ↓ Scompare il contributo del tempo
 ↓ i ∈ S.
 ↓ $\gamma < 1$ e' un operatore di contrazione ⇒ ∃! punto fisso

If the immediate contribution is stochastic, possibly because it depends on the next state, we may denote it as $h(i, a, j)$ and rewrite Eq. (3) as

Ottimizzazione fuori
Valore atteso dentro

$$V(i) = \underset{a \in \mathcal{A}(i)}{\text{opt}} \sum_{j \in \mathcal{S}} \pi(i, a, j) \left\{ h(i, a, j) + \gamma V(j) \right\}, \quad i \in \mathcal{S}. \quad (4)$$

↓ Non e' detto che funzioni sempre.
 ↓ Meglio la scrittura $h(i, a, j), \gamma, V(j)$
 ↓ Contributo immediato stocastico che dipende dallo stato in cui vado a finire dopo l'istante di tempo
 ↓ realizzazione dei fattori di rischio
 ↓ Contro Es

In principle, we may write

Contro Es

- Se consideriamo il livello di magazzino h non dipenderà solo da i, a, j . Se vado a finire in $I_c = 0$ non sai come ci sei arrivato, cioè $I_c = dt$ o $I_c < dt$ ma in queste opzioni portano contributi immediati differenti

but this may be not practical, either because N is finite but huge, or because transition probabilities are not known.

$$f(i, a) = \sum_{j \in \mathcal{S}} \pi(i, a, j) h(i, a, j),$$

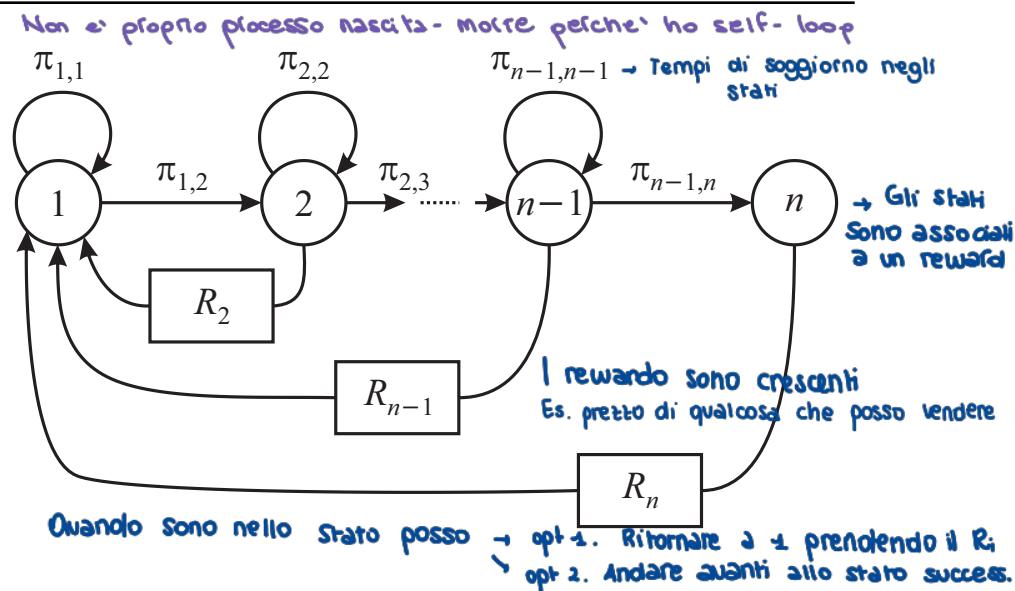
↓ Do un voto all'azione e lo stato finale in cui sei finito
 ↓ Devi conoscere questi valori
 ↓ Non e' detto che le conosca, le posso stimare

via di mezzo tra apprendimento
supervisionato oppure no
 ↑
 (in reinforcement learning)
 non li conosco

Example recurrent optimal stopping

To illustrate finite MDPs and controlled Markov chains let us consider a simple case, where states $k = 1, 2, \dots, n$ are arranged in a sort of linear chain.

The state can only move to the right, but it can be reset, resulting in the collection of a reward and a system restart.



When we are at state 1, we will stay there with probability π_{11} , or we will move to the next state 2 with probability π_{12} . In state 1, there is no decision to make.

In each “interior” state $k = 2, \dots, n - 1$, we have a choice between two actions: wait and reset. → Aspetto e esercito subito?

- If we wait and do nothing, there is no reward and we remain in state k with probability π_{kk} , or we move to the next state in the chain with probability $\pi_{k,k+1}$.
- Alternatively, we may exercise an option to collect an immediate reward R_k and reset the state to 1.

Since the transition is deterministic when we choose the reset action, we streamline the notation and use π_{ij} as a shorthand for $\pi(i, \text{wait}, j)$, $i, j \in \mathcal{S}$.

With reference to Eq. (3), the immediate contributions are

$$f(k, \text{wait}) = 0, \quad f(k, \text{reset}) = R_k.$$

Aspetti / Speri di finire in uno stato migliore

We assume that the rewards are increasing along the chain, i.e., $R_k < R_{k+1}$, so there is a tradeoff between collecting the immediate reward and waiting for better opportunities. If we are patient enough and reach the last state n , we may immediately earn the largest reward R_n and move back to state 1.

This problem bears a definite similarity with optimal stopping problems. A possible interpretation is that we own an asset, whose price process is stochastic, and we have to choose when it is optimal to sell the asset. In general, the price need not be monotonically increasing.

Il problema e' stazionario \Rightarrow queste sono equazioni
 The DP equations to find the value $V(i)$ of each state $i \in \mathcal{S}$ are as follows:

$$V(1) = \gamma\pi_{1,1}V(1) + \gamma\pi_{1,2}V(2)$$

Fattore di sconto
Rimango nello stato 1

$$V(k) = \max \left\{ R_k + \gamma V(1), \underbrace{\gamma\pi_{k,k}V(k) + \gamma\pi_{k,k+1}V(k+1)}_{\text{Opzioni dei nodi intermedi}} \right\}, \quad k = 2, \dots, n-1$$

reset
Value function associata alla politica ottima
Vado avanti allo stato 2

$$V(n) = R_n + \gamma V(1). \rightarrow \text{Torni indietro prendendo il Reward immediato}$$

wait
Risoluzioni di questo sistema di equazioni lineari a tratti (causa max) fattibile con il metodo iterativo value iteration

Policy evaluation and Q -factors

The previous example might suggest that finite MDPs are a rather simple business to deal with.

→ perché spazio degli stati finito
The value function is just a vector and, in the case of a finite horizon, we have an explicit relationship between value functions at different stages. If the set of feasible actions is not too large, the optimization step is trivially solved by enumeration.

In the infinite-horizon case, we have to find state values by solving a system of equations, which are not linear (piecewise linear, in fact). However, they can be solved by rather simple iterative methods. → *uno potrebbe avere l'illusione che sia semplice: No!*

1. Spazio degli stati
2. Trovare $\pi(i, a, j)$
3. Conoscere il modello e/o la dinamica

However, the MDP business may not be so simple, because of the sheer size of the state space. Furthermore, even finding the whole set of transition probabilities $\pi(i, a, j)$ may be quite difficult, if not impossible, when the curse of modeling strikes.

One ingredient to circumvent these difficulties is •Monte Carlo sampling. Another ingredient is to rewrite the DP recursion in a different way, based on Q -factors. Let us do this for the infinite horizon case.

collegato allo stato e all'azione: misura quanto e' buona l'azione
Informally, a Q -factor $Q(i, a)$ measures the value of taking action a when in state i . More precisely, this should be defined with reference to the policy that will be applied in the next steps. → *Ma al prossimo step? Voglio capire quanto la politica e' buona → vogliamo fare learning*

Let us consider a feasible stationary policy μ , mapping each state i into action $a = \mu(i) \in \mathcal{A}(i)$. The state value function when following policy μ can be found by (solving a system of linear equations):

Fissando μ otteniamo un'equazione per ogni stato, per ottenere sol bisogna risolvere il sistema lineare (affine)

La politica mappa lo stato in un'azione

Nell'immediato seguo μ

Value associato a una politica (non necessariamente ottima)

$$V_\mu(i) = f(i, \mu(i)) + \gamma \sum_{j \in S} \pi(i, \mu(i), j) \cdot V_\mu(j), \quad i \in S. \quad (5)$$

We will discuss Eq. (5) in more depth later.

For now, it suffices to say that it is fundamental for numerical methods, based on policy iteration, in which we first assess the value of a candidate stationary policy μ , and then we try to improve it.

Ragionano sul cambiare la politica stazionaria → Prendi una politica, la valuti, capisci quanto è buona, poi la valuti

Also note the difference with respect to Eq. (3), which involves the optimal choice of action a . Here, there is no optimization involved, as the action is prescribed by the selected policy μ .

We define the Q -factor for the stationary policy μ as the following mapping $S \times \mathcal{A} \rightarrow \mathbb{R}$:

Politica qualsiasi

politica

Qui invece seguo la politica μ

Sono in i e scelgo a senza seguire per forza μ

Probabilità di partire da i seguire a e arrivare a j

$$Q_\mu(i, a) \doteq f(i, a) + \gamma \sum_{j \in S} \pi(i, a, j) \cdot V_\mu(j). \quad (6)$$

The idea is to apply action a at the current state i , and then follow policy μ .

We may also define the optimal Q -factors by plugging the optimal value function into Eq. (6):

Politica ottima

qui uso la/le politiche ottime

$$Q(i, a) = f(i, a) + \gamma \sum_{j \in S} \pi(i, a, j) V(j), \quad i \in S, a \in \mathcal{A}(i). \quad (7)$$

Then,

Se avessi tutti i a -factor, la politica ottimale
Sarebbe questa $V(j) \equiv \underset{a \in A(j)}{\text{opt}} Q(j, a), \quad j \in S.$

Hence, we may rewrite the DP recursion in terms of Q -factors:

$$Q(i, a) = f(i, a) + \gamma \sum_{j \in S} \pi(i, a, j) \left[\underset{\tilde{a} \in A(j)}{\text{opt}} Q(j, \tilde{a}) \right], \quad i \in S, a \in A(i). \quad (8)$$

• Non facilmente capibili

↓

Tiriamo fuori il corrispettivo statistico del prob. probabilistico

If we compare Eqs. (3) and (8), we may notice that there are both good and bad news:

- The bad news is that, instead of a state value function $V(i)$, now we have state-action value functions $Q(i, a)$. Unless we are dealing with a small sized problem, it seems that we have just made the curse of dimensionality even worse.
↳ Sembra di aver peggiorato
- The good news is that now we have swapped expectation and optimization.
↳ Ma! Buona idea di scambiare $\mathbb{E} = \text{opt}$

We will appreciate the advantage of swapping expectation and optimization later, especially when dealing with continuous problems: it is often easier to solve many simple (deterministic) optimization problems than a single large (stochastic) one.

Furthermore, we may learn the Q -factors by statistical sampling, which is useful in both large-scale problems and problems with an unknown underlying dynamics. This paves the way to model-free DP.

Last, but not least, when dimensionality precludes finding the factors $Q(i, a)$ for all state-action pairs, we may devise a compact representation based on an approximation architecture, which may range from a relatively straightforward linear regression to deep neural networks. → Ogni problema richiede il suo strumento

Different shades of stochastic DP

Let us consider once more the basic recursive DP equation, as given in Eq. (1):

$$V_t(s_t) = \underset{x_t \in \mathcal{X}(s_t)}{\text{opt}} \left\{ f_t(s_t, x_t) + \gamma \mathbb{E}[V_{t+1}(s_{t+1}) \mid s_t, x_t] \right\},$$

Non è il modello più generale, né computazionalmente più efficiente \Rightarrow dipende dal problema!
 1. Integrale in molte dim.
 Equazione a orizzonte infinito Rognosi: 2. Hai i valori di π !
 Per trovare V_t bisogna risolvere vari problemi stocastici
 Le distribuzioni potrebbero essere nel continuo

The underlying assumptions are:

1. We first observe the state s_t at time instant t .
2. Then, we make a feasible decision $x_t \in \mathcal{X}(s_t)$.
3. Then, we observe an immediate contribution $f_t(s_t, x_t)$.
4. Finally, we move to a new state s_{t+1} , which depends on realized risk factors ξ_{t+1} during the subsequent time interval, according to a probability distribution that might depend on the current state s_t and the selected decision x_t .

In order to find the value $V_t(\cdot)$ of each state s_t , based on knowledge of $V_{t+1}(\cdot)$, we should solve a stochastic optimization problem that may involve a quite challenging expectation. However, in the case of an MDP, we have seen that Q -factors allow us to swap optimization and expectation. Actually, we may have to swap them in order to reflect the *information structure*, as shown in the following example.

Incertezza bassa nell'immediato, alta a lungo termine $\Rightarrow d_t$ e' nota al tempo $t-1$

→ Hai idea delle immediate domande future

Example: Lot-sizing with limited lookahead

Per ogni istante di tempo osservo I_t → decido quanto ordinare
il magazzino diventa $I_t^* = I_t + x_t$, osservo la domanda d_{t+1}
Quindi $I_{t+1} = I_t^* - d_{t+1}$

stato pre-decisione

stato post-decisione

Under demand uncertainty, we should wonder about the precise nature of the demand process $(d_t)_{\{t \geq 1\}}$. → Come è fatto questo processo stocastico?

In che contesto sono?

Il cliente è un'altra azienda, quindi ci sono ordini formali

→ Magazzino svuotato dal consumatore finale, non ci sono ordini
In a B2C (business-to-consumer) setting, like a retail store, we may not have any advance information about demand in the next time interval. However, in a B2B (business-to-business) setting, we probably receive formal orders from customers over time. In the limit, we may have perfect information about demand in the next time interval, even though uncertainty may be relevant for the not-so-distant future.

If, at time instant t , we have perfect knowledge about demand d_{t+1} , the immediate cost contribution is actually deterministic:

Anche se la domanda è deterministica non è detto che la vogli sempre soddisfare (Φ)

$$f_t(I_t, x_t, d_{t+1}) = hI_t + cx_t + \phi \cdot \delta(x_t) + q \cdot \max_{\text{costo fisso ordinazione}} \{0, d_{t+1} - I_t - x_t\},$$

Il contributo immediato è deterministico: conosco d_{t+1}

Non è detto che in ogni periodo si ordini (costo fisso ϕ)

Includiamo una potenziale vendita persa da minimizzare

involving an inventory holding charge h , variable and fixed ordering costs c and ϕ , and a lost sales penalty q . Note that if there is a large fixed ordering charge, we might prefer to leave a small amount of demand unsatisfied.

In order to write the DP recursion, we should pay attention to the precise relationships among information flow, decisions, and system dynamics. Given the current inventory state I_t , first we observe demand, then we make a decision. Hence, the DP recursion in this case reads:

$$V_t(I_t) = E_t \left[\min_{x_t \geq 0} \left\{ f_t(I_t, x_t, d_{t+1}) + V_{t+1}(I_{t+1}) \right\} \right].$$

Mi serve il valore condizionato rispetto alle info che ho in questo istante

Ora non so che domande future osserverò

→ Ora vedo l'effetto del look ahead limitato: ad ogni step ho un problema deterministico e incerto 13 solo sugli step successivi

Costruisco il problema in maniera diversa introducendo la variabile post-decisione, perche' scrivendola con opt all'interno, questa e' deterministica e il IE e' trouvable con approcci statistici

Post-decision state variables

The previous example shows a case in which the form of the DP recursion is

In funzione della dinamica del sistema
il valore atteso esso dentro o fuori

$$V_t(s_t) = E_t \left[\underset{x_t \in \mathcal{X}(s_t)}{\text{opt}} \left\{ f_t(x_t, s_t) + \gamma V_{t+1}(s_{t+1}) \right\} \right]. \quad (9)$$

The potential advantages of this form are:

- The inside optimization problem is deterministic.
- The outside expectation may be estimated by statistical sampling.

Since these are relevant advantages, we may try to rephrase the standard DP recursion in this form, even though it may not be the natural one.

Sometimes, this can be accomplished by rewriting the system dynamics based on post-decision state variables.

In the familiar view of system dynamics, we consider transitions from state s_t to state s_{t+1} , depending on the selected decision x_t and the realized risk factor ξ_{t+1} . Sometimes, it is convenient to introduce an intermediate state, observed after the decision x_t is made, but before the risk factor ξ_{t+1} is realized. Such a state is referred to as **post-decision state**, and it may be denoted by s_t^x .

Idea: vedere le transizioni in due step

Example. In a lot-sizing problem, we may introduce a post-decision state, the inventory level I_t^x after making the ordering decision at time instant t , but before observing and meeting demand during the time interval $t + 1$:

Spezziamo di dividere la transizione in due

$$I_t^x = I_t + x_t. \quad \text{Magazzino pre-decisione}$$

Hence, we break the transition to the next state into two steps, the second one being

$$I_{t+1} = I_t^x - d_{t+1}. \quad \text{Magazzino post-decisione}$$

The introduction of post-decision states changes the dynamics from the standard sequence

$$(s_0, x_0, \xi_1; s_1, x_1, \xi_2; \dots; s_{t-1}, x_{t-1}, \xi_t; \dots)$$

to

Stato azione
Fattori esogeni di rischio
Stato intermedio → post-decision state

$$(s_0, x_0, s_0^x, \xi_1; s_1, x_1, s_1^x, \xi_2; \dots; s_{t-1}, x_{t-1}, s_{t-1}^x, \xi_t; \dots)$$

In other words, we split the transition equation $s_{t+1} = g_{t+1}(s_t, x_t, \xi_{t+1})$ into two steps:

Equazioni di transizione

$$s_t^x = g_t^1(s_t, x_t), \quad \rightarrow \text{Transizione iniziale} \quad (10)$$

$$s_{t+1} = g_{t+1}^2(s_t^x, \xi_{t+1}). \quad \rightarrow \text{Nuova transizione} \quad (11)$$

This is also reflected in terms of value functions, as we may introduce the **value of the post-decision state** at time t , $V_t^x(s_t^x)$, where superscripts x point out that these quantities refer to post-decision states. Usiamo la value function con stato post decisione

The relationship between the standard value function $V_t(s_t)$ and $V_t^x(s_t^x)$ can be expressed as

La value function e' interpretabile come il valore atteso

$$V_t^x(s_t^x) = \mathbb{E}[V_{t+1}(s_{t+1}) | s_t^x]. \quad (12)$$

l'effetto della decisione presa e' qui dentro

Note that, in moving from s_t^x to s_{t+1} , no decision is involved. Then, using Eq. (12), we may rewrite the standard DP recursion as a deterministic optimization problem,

$$V_t(s_t) = \underset{x_t \in \mathcal{X}(s_t)}{\text{opt}} \left\{ f_t(s_t, x_t) + \gamma V_t^x(s_t^x) \right\}, \quad \rightarrow \text{Approssiamo il IE con una regressione lineare} \quad (13)$$

where the post-decision state s_t^x is given by the first transition equation (10). If we move one step backward in time to instant $t - 1$ and write Eq. (12) for state s_{t-1}^x , we have

$$V_{t-1}^x(s_{t-1}^x) = \mathbb{E}[V_t(s_t) | s_{t-1}^x]. \quad \rightarrow \text{Scriviamo la Value function all'istante di tempo precedente (t-1)}$$

Then, by plugging Eq. (13) into this relationship, we obtain

$$V_{t-1}^x(s_{t-1}^x) = \mathbb{E}[V_t(s_t) | s_{t-1}^x]$$

*con il post decisione
Siamo riusciti a portare
il valore atteso fuori*

$$= \mathbb{E} \left\{ \underset{x_t \in \mathcal{X}(s_t)}{\text{opt}} \left[f_t(s_t, x_t) + \gamma V_t^x(s_t^x) \right] \mid s_{t-1}^x \right\}. \quad (14)$$

Leaving the inconsequential backshift in time aside, we find again a DP recursion featuring a swap between expectation and optimization.

The DP recursion in terms of Q -factors, actually, fits naturally into this framework, since the state-action pair (i, a) may be regarded as a post-decision state, before observing the random transition to the new state.

State augmentation in inventory management

on-order: è il magazzino che tiene conto della merce ordinata, ma non ancora arrivata

Quello che hai a portata

When dealing with inventory management, on-hand inventory looks like the natural state variable and the obvious basis for any ordering decision.

On the contrary, ordering decisions should not be based on on-hand inventory, but on available inventory, which accounts for inventory on-order (items ordered from the supplier, but not yet delivered) and backlog.

The standard state equation

$$I_{t+1} = I_t + x_t - d_{t+1}$$

assumes that what is ordered at time instant t is immediately available to satisfy demand during the following time interval $t+1$.

If the delivery lead time is an integer number $LT \geq 1$ of time intervals, we need to introduce state variables keeping track of what was ordered in the past and is still in the transportation pipeline.

Let us denote these state variables by $z_{t,\tau}$, the amount that will be delivered τ time intervals after the current time t , where $\tau = 0, 1, 2, \dots, LT-1$. Hence, $z_{t,0}$ represents what is immediately available at the current time instant t , and it is involved in the state transition equation for on-hand inventory:

$$I_{t+1} = I_t + z_{t,0} - d_{t+1},$$

if we disregard demand uncertainty.

What we order at time instant t , represented by decision variable x_t , will be available LT time intervals after t . Hence, at the next time instant $t + 1$, the amount x_t will be LT – 1 time intervals from delivery. We may therefore relate the decision x_t to the additional state variable corresponding to $\tau = LT - 1$ as follows:

$$z_{t+1,LT-1} = x_t.$$

↗ Variabile di controllo
↙ Variabile di stato

The general transition equation for the additional state variables $z_{t,\tau}$, for $\tau < LT - 1$, boils down to a simple time shift:

$$z_{t+1,\tau} = z_{t,\tau+1}, \quad \tau = 0, 1, \dots, LT - 2. \quad (15)$$

↗ Oggetti deperibili: scadenze, rotture, ...

By a similar token, when dealing with perishable items we need to describe inventory at each age, by introducing an array of state variables $I_{t,\tau}$ that represent the amount of on-hand inventory at time t with an age of τ time periods.

The state transitions are not quite trivial to write down (see the book), and depend on the kind of inventory issuing, **FIFO** (first-in-first-out) or **LIFO** (last-in-first-out).

Al tempo t quanti oggetti hanno eta' pari a 2 dopo un po' l'oggetto non e' più vendibile

Revenue management

Ricavo, non profitto → il costo lo consideriamo "assorbito" o trascurabile

Revenue management consists of a series of models and techniques to maximize the revenue obtained by selling perishable resources. The idea was introduced (under the name of yield management) within the airline industry.

La variabile decisionale è
il prezzo

Trategia basata sulla quantità

There are two basic approaches to revenue management: quantity-based and price-based. In the first case, resource availability is restricted according to some policy in order to maximize revenue. In the second one, prices are adjusted dynamically, and we talk of dynamic pricing.

We outline simple DP models for quantity-based revenue management and assume that the marginal cost of each unit of resource is zero or negligible, so that profit maximization boils down to revenue maximization.

Ricorda

Regola di Littlewood:

+ riguardiamo per
programmazione dinamica

We consider C units of a single resource (say, seats on an aircraft), which must be allocated to n fare classes, indexed by $j = 1, 2, \dots, n$. Units in class j are sold at price p_j , where

diverse classi di tariffa

Prima classe

$$p_1 > p_2 > \dots > p_n$$

n -esima classe

Domande: • Il mercato è segmentato?

• Come la distribuzione delle domande?

Thus, class 1 is the first class, from which the highest revenue is earned.

We assume that the seats allocated to different classes are actually identical, but they are bundled with ancillaries (like cancellation rights, weekend restrictions, on-board meals, etc.) to offer different packages.

Vendi il posto e alcuni servizi / Pasti a bordo / Bagaglio aggiuntivo

The price of each class is fixed, but we control the available inventory by restricting the seats available for each class.

✓ E' il fattore di rischio per ogni classe

Demand D_j for each class is random, but it may be realized according to different patterns, giving rise to different problem formulations.

The two essential features that we consider here are:

✓ Ogni passeggero «ha in testa» una singola classe → non devo costruire un modello di scelta

- **Customer behavior.** In a perfectly segmented market, any passenger is willing to buy only one kind of class. Hence, we do not need to model passengers' preferences. On the contrary, if we want to account for preferences, we must define a choice model. → Aspetto più realistico

✓ Le domande arrivano tutte insieme oppure lungo un periodo in modo sequenziale

- **Timing of demand.** It would be nice if demand occurred sequentially, class 1 first and then down to class n . The worst case, in some sense, is when low-budget customers arrive first.

A model where there exist disjoint time intervals, during which exactly one class is in demand, is called "static." The term "dynamic" is reserved to models where customer requests for different classes are interleaved over time.

Arrivi per classi
(prima tutta di una, poi la successiva, etc...)

Arrivi mescolati

We can express the decision policy in terms of maximum amount of seats available per class (protection levels). Alternatively, we could use bid-prices (i.e., the minimum price at which we are willing to sell a seat).

Questo è ciò che scelgo io per proteggere le classi più costose

→ Realizzazione sequenziale di variabili casuali indipendenti per il mercato perfettamente segmentato

Static model with perfect demand segmentation

Demand for classes are assumed independent and occur sequentially, D_n first. No customer is willing to switch to another class.

Therefore, we may associate decision stages with classes, indexed in decreasing order of the index $j = n, n - 1, \dots, 2, 1$.

↳ indici di stadio (decisionali)

The natural state variable is s_j , the integer number representing the residual capacity at the beginning of stage j ; (the initial state is $s_n = C$, the number of available seats on the aircraft.)

Quanti posti hai sull'aereo prima di osservare la domanda

In order to decide at each stage j how much of the remaining capacity is offered to class j , we need to find a sequence of value functions $V_j(s_j)$. → Qual è il valore di avere allo stadio j
Sì posti disponibili

The overall objective is to find the maximum expected revenue that we may collect by selling the C available seats to the n classes, denoted by $V_n(C)$, and the boundary condition is

$$V_0(s_0) = 0, \quad s_0 = 0, 1, \dots, C.$$

I posti vuoti sono stati «sprecati»

↳ Valore atteso del ricavo seguendo la strategia all'ottimo

We do not consider groups or families and, given independence, there is no demand learning effect along the way.

A rather surprising finding is that we may build the DP recursion by assuming the following event sequence, for each class j :

Assumiamo questa come sequenza della dinamica

1. First, we observe demand D_j for class j .
2. Then, we decide how many requests for class j to accept, a decision represented by the integer number x_j .
3. Then, we collect revenue $p_j x_j$ and proceed to stage $j - 1$, with the updated state variable $s_{j-1} = s_j - x_j$. → Equazione di transizione di stato

This does not seem plausible, as we are supposed to make decisions before observing demand. However, we shall prove that the decision x_j does not rely on the probability distribution of D_j .

Non dipende dalla distribuzione di probabilità di

The intuition is that we do not really need to declare x_j beforehand. We may observe each individual request for class j sequentially, and each time we decide whether to accept it or not. If we reject a request, then we declare class j closed.

Formally, the DP recursion takes the swapped form that we introduced in Eq. (9):

Regola di Littlewood
quantile della distribuzione
distributione di probabilità della domanda
 $y_i^* = F_i^{-1}(1 - R_i)$
Livello di protezione della classe 1

NB Non entra la distribuzione della classe 2
ENSR : euristica per più classi

$$V_j(s_j) = E \left[\max_{0 \leq x_j \leq \min\{s_j, D_j\}} (p_j x_j + V_{j-1}(s_j - x_j)) \right]$$

I j sono decrescenti : da n a 1
Tiro fuori l'equazione come se avessi
Di prima rispetto alla decisione
→ Anche qui IE fuori e max dentro

Quanti posti assegno alla classe j
E' fuori perché scelgo y_j conoscendo D_j
Nei prendere la decisione conosco D_j
Non decido prima quanti posti assegnare, ma mano mano
che arrivano i clienti

Introduciamo uno shift dell'indice

As it turns out, it may be convenient to introduce an inconsequential index shift in j and omit the stage subscript in decision variable x and state s to ease notation:

$$V_{j+1}(s) = E \left[\max_{0 \leq x \leq \min\{s, D_{j+1}\}} (p_{j+1}x + V_j(s-x)) \right]. \quad (16)$$

The decision x at stage $j+1$ is naturally bounded by the residual capacity s and by the realized demand D_{j+1} .

In order to analyze the structure of the optimal policy and justify the above form, it is convenient to introduce the expected marginal value of capacity

$$\Delta V_j(s) \doteq V_j(s) - V_j(s-1),$$

Se $p_j > \Delta V_j(s)$ mi conviene vendere il biglietto
 perché il ricavo ora vale di più di vita vendita funz.
 E' il valore di avere s posti per la classe
 j meno il valore di averne $s-1$
 Si lega alla perdita di valore
 E' il costo opportunita' di avere un posto in meno
 cioè quanto ci perdo nel futuro a vendere ora
 un biglietto della classe j

for each stage j and residual capacity s .

Intuitively, the marginal value of capacity measures, for a given residual capacity s , the opportunity cost of an aircraft seat, i.e., how much revenue we are expected to lose if we give up a seat.

Then, we may rewrite Eq. (16) as

$$V_{j+1}(s) = V_j(s) + E \left[\max_{0 \leq x \leq \min\{s, D_{j+1}\}} \left\{ \sum_{z=1}^x (p_{j+1} - \Delta V_j(s+1-z)) \right\} \right] \quad (17)$$

quando la differenza si annulla, hai trovato quanti biglietti vendere
 mi conviene vendere
 non mi conviene vendere
 Questa differenza e' positiva o negativa?
 e' decrescente rispetto a z
 => All'aumentare di 2, la differenza parte positiva
 e poi diventa negativa
 e' uguale alla moltiplicazione (!)
 Conta i posti, da z a z , che si vendono alla classe j importante
 Amico a un valore pari a s (col cambio di segno per la differenza)
 formulazione finale

Notiamo che risolviamo il max per ogni valore di D_j possibile e poi faccio il IE pesando il valore ottenuto con le probabilità di D_j
 Mi fermo Sbarco sul min

This rewriting is based on a standard trick of the trade, i.e., a telescoping sum:

sfruttiamo la somma telescopica

$$\begin{aligned}
 V_j(s-x) &= V_j(s) - [V_j(s) - V_j(s-x)] \rightarrow \text{aggiungo e tolgo } V_j(s) \\
 &= V_j(s) - [V_j(s) \pm V_j(s-1) \pm \cdots \pm V_j(s+1-x) - V_j(s-x)] \\
 &\quad \downarrow \text{Significa } + V_j(s-1) - V_j(s-1) \\
 &= V_j(s) - \sum_{z=1}^x [V_j(s+1-z) - V_j(s-z)] \rightarrow \text{Esprimo come incrementi successivi} \\
 &\quad \rightarrow \text{Tutti gli altri fatti si elidono (somma telescopica!)} \\
 &= V_j(s) - \sum_{z=1}^x \Delta V_j(s+1-z), \rightarrow \text{Somma dipendente da } z \\
 &\quad \downarrow \text{Serie di differente scritte sopra} \\
 &\quad \text{Non dipende da } z \text{ quindi lo posso portare fuori}
 \end{aligned}$$

where we use the shorthand $\pm V \equiv +V - V$. By plugging the last expression into Eq. (16) we easily obtain Eq. (17).

The following facts about expected marginal values can be proved:

Il valore marginale è decrescente diversamente in base a quanti potrò ho => proprietà di monotonia

Lasciare un posto quando ne ho tanti è meno «doloroso» rispetto a quando ne ho pochi

$$\Delta V_j(s+1) \leq \Delta V_j(s), \quad \Delta V_{j+1}(s) \geq \Delta V_j(s), \quad \forall s, j.$$

Proprietà di MONOTONIA

The idea is that marginal values are decreasing with available capacity and that the value of capacity is larger when more stages are remaining.

As a consequence, the terms

$$p_{j+1} - \Delta V_j(s+1-z) \rightarrow \text{termine decrescente in } z$$

in the sum of Eq. (17) are decreasing with respect to z .

Cumulare termini nella somma = Continuo a soddisfare le richieste

To find the optimal solution, we should keep increasing x , i.e., cumulating terms in the sum and accepting ticket requests for class $j+1$, until we find the first negative term or the upper bound of the sum, $\min\{s, D_{j+1}\}$, is reached.

We should keep accepting offers while the revenue p_{j+1} is larger than the opportunity cost measured by the marginal value of a seat ΔV_j for the next classes.

Decido il livello max di posti che assegnerò a una classe, non quanti effettivamente ne renderò

This defines an optimal nested protection level,

La politica e' annidata

E' il livello di protezione che soddisfa la condizione, così definisco la struttura della politica ottima

$$y_j^* = \max\{y : p_{j+1} < \Delta V_j(y)\}, \quad j = 1, \dots, n-1. \quad (18)$$

y Ammontare dei posti da j fino a j

This is the amount of seats reserved to classes $j, j-1, \dots, 1$; the term "nested" is due to the fact that capacity is allocated to all classes higher than j included, rather than to a single class.

Equation (18) says that we should increase the protection level y_j until we find that the revenue p_{j+1} from the previous class $j+1$ compensates the opportunity cost of a seat.

The optimal protection level for these next classes constrains the number of tickets we may sell for class $j+1$, and we may express the optimal decision at stage $j+1$ as

$$x_{j+1}^*(s_{j+1}, D_{j+1}) = \min\{(s_{j+1} - y_j^*)^+, D_{j+1}\}. \quad \begin{array}{l} \text{massimo che posso offrire a questa classe} \\ \text{potrebbe essere negativa (anche se non dovrebbe)} \\ \rightarrow \text{E' la decisione finale data in funzione della domanda} \end{array}$$

This means that the maximum number of seats sold to class $j+1$ is the minimum between the observed demand and the excess residual capacity with respect to the protection level for the higher classes.

Indeed, we do not need to know D_{j+1} in advance, since we may keep satisfying demand for class $j+1$ until either the stage ends and we move on to the next class, or the protection level is reached, or we just run out of seats. This justifies the swap of expectation and minimization in Eq. (17). *→ Facendo questi conti ottengo gli y_j^* e ottengo così una politica che deve essere flessibile rispetto alle richieste dinamiche che mi arrivano*

Cosa succede se rilassiamo l'assunzione degli arrivi disgiunti per le varie classi, lasciamo invece il mercato segmentato
Assumiamo sempre che i passeggeri siano indipendenti

Dynamic model with perfect demand segmentation

Let us relax the assumption that demand for classes occur sequentially over disjoint time intervals. However, we still assume a rigid market segmentation.

Dobbiamo mettere il tempo, che sarà discretizzato → osservo al più un arrivo in ogni intervallo di tempo → processo di Poisson

Time is involved in the DP recursion, but a simple model is obtained if we assume that time intervals are small enough that we may observe at most one arrival per time interval.

Assumiamo t come indice intero

Let us denote by $\lambda_j(t)$ the probability of an arrival for class j during time interval t , $t = 1, \dots, T$. Since we assume a perfectly segmented market, these probabilities refer to independent events and must satisfy the consistency constraint:

$$\sum_{j=1}^n \lambda_j(t) \leq 1, \quad \forall t.$$

Perche' potrebbero non arrivare richieste del S_k di tempo
↪ 0 arriva un evento oppure non ne arrivano

The DP recursion in this case aims at finding value functions $V_t(s)$, where s is residual capacity, subject to the boundary conditions

$$V_t(0) = 0, \quad t = 1, \dots, T, \quad \rightarrow \text{Se non ha più posti disponibili, il valore è 0}$$
$$V_{T+1}(s) = 0, \quad s = 0, 1, \dots, C, \quad \rightarrow \text{Dopo il decollo i posti non utilizzati non valgono più nulla}$$

where $T + 1$ denotes the end of the time horizon (when the aircraft takes off).

Ricordati che il mercato e' perfettamente segmentato
 Se vendi il biglietto ottieni questo, ma non hai ancora deciso di vendere il biglietto

We may denote by $R(t)$ the available revenue at time t , which is a random variable taking value p_j when an arrival of class j occurs, 0 otherwise.

Whenever there is an arrival, we have to decide whether we accept the request or not, a decision that may be denoted by the binary decision variable x .

As in the static model, we do not need to make a decision in advance, since we just need to react to a request. Hence, the DP recursion takes again the "swapped" form

Ragioniamo una richiesta per volta

$$V_t(s) = E \left\{ \max_{x \in \{0,1\}} \left[R(t)x + V_{t+1}(s-x) \right] \right\}.$$

La ragione e' che il mercato e' perfettamente segmentato
 ↳ Accetto oppure no la richiesta

→ Posso osservare la richiesta e dopo decidere se accettare o rifiutare

Using expected marginal values of capacity again, a policy in terms of protection levels may be devised. In this case, however, protection levels may be time-varying.

E' il modello che rilascia tutte le assunzioni precedenti
Le richieste sono mescolate (il modello e' dinamico) e il cliente sceglie la classe (chi arriva puo' scegliere quale biglietto prendere)

Dynamic model with customer choice

→ discrete choice: se ho a disposizione più prodotti quali scelgo

If we relax the assumption of perfectly segmented markets, then we must define a model of customer choice.
Assumiamo che qualcuno ci dia questo modello di scelta
La variabile di controllo e' al tempo t quante classi tariffe si stanno offrendo

One way of doing so is to partition the market into segments and to estimate the fraction of passengers belonging to each segment, as well as the probability that a class will be purchased by a passenger of each segment, if she is offered that class.

↪ A ogni segmento mettiamo la probabilità che compra un oggetto di una determinata classe

The control decision, at each time instant t , is the subset of classes that we offer.

Let \mathcal{N} be the set of available classes, indexed by j and associated with revenue p_j . It is convenient to introduce $p_0 = 0$, the zero revenue earned when there is a passenger willing to fly but, given the subset of classes currently offered, she decides not to purchase any ticket.

Formally, a potential passenger request arrives at time t with probability λ and, given the offered subset of classes $\mathcal{S}_t \subseteq \mathcal{N} = \{1, \dots, n\}$, she will make a choice.

↪ Probabilità che un cliente sceglia la classe j in relazione a quello che si sta offrendo in questo momento
Let $P_j(\mathcal{S}_t)$ be the probability that the customer chooses class j , as a function of the offered set, where we include the probability $P_0(\mathcal{S}_t)$ of no-purchase.

These probabilities may be estimated, based on the aforementioned market segmentation model, and they are subject to the natural constraints

Non e' per nulla ovvio come ottenere queste info parlando dall'analisi di mercato

$$P_j(\mathcal{S}) \geq 0,$$

$$\sum_{j \in \mathcal{S}} P_j(\mathcal{S}) + P_0(\mathcal{S}) = 1,$$

$$\mathcal{S} \subseteq \mathcal{N}, j \in \mathcal{S} \cup \{0\}$$

$$\mathcal{S} \subseteq \mathcal{N}. \rightarrow \text{Per ogni sottinsieme di tutte le classi le probabilità devono sommare a 1}$$

In this model, the decision variable is actually a set, i.e., the subset \mathcal{S}_t of classes offered at time t .

The decision-dependent purchase probabilities at time t are

$$\begin{aligned} & \text{Probabilità di avere un arrivo} \\ & \lambda P_j(\mathcal{S}_t), \quad j = 1, \dots, n \\ & (1 - \lambda) + \lambda P_0(\mathcal{S}_t), \quad \text{Se il cliente non compra nulla} \quad j = 0. \\ & \quad \text{Se non c'è arrivo nessuno} \end{aligned}$$

In the last case, we may have a no purchase either because no one came, or because the passenger did not like the offered assortment of choices.

The DP recursion is

$$V_t(s) = \max_{\mathcal{S}_t \subseteq \mathcal{N}} \left\{ \sum_{j \in \mathcal{S}_t} \lambda P_j(\mathcal{S}_t) (p_j + V_{t+1}(s-1)) + (\lambda P_0(\mathcal{S}_t) + 1 - \lambda) V_{t+1}(s) \right\}.$$

Il massimo è fuori dal IE perché il mercato non è segmentato e il cliente quando arriva può decidere quale prodotto vuole comprare => l'omogeneizzazione è fuori dal valore atteso
se vendo al periodo successivo ho un posto in meno

la mia decisione sta nella tipologia di offerta che faccio => prima offro, poi il cliente decide cosa accadrà
non vendo → ho ricavo 0 → ma ho un posto in più all'incontro successivo

In this case, we must lay our cards down on the table before the passenger's decision. As a consequence, the DP recursion is in the usual form, featuring the maximization of an expectation.