

# Reinforcement Learning and Artificial Intelligence in the context of revenue management

Freie wissenschaftliche Arbeit  
zur Erlangung des akademischen Grades  
“Master of Science”  
Studiengang: Finanz- und Informationsmanagement

an der  
Wirtschaftswissenschaftlichen Fakultät  
der Universität Augsburg

– Lehrstuhl für Analytics & Optimization –

Eingereicht bei:	Prof. Dr. Robert Klein
Betreuer:	tbd
Vorgelegt von:	Stefan Glogger
Adresse:	tbd
Matrikel-Nr.:	tbd
E-Mail:	stefan.glogger@student.uni-augsburg.de
Datum:	July 9, 2019

# Contents

<b>1</b>	<b>Aufzuschreiben</b>	<b>3</b>
<b>2</b>	<b>Problem and Approaches to Solution</b>	<b>4</b>
2.1	The problem . . . . .	4
2.1.1	Dynamic Programming . . . . .	4
2.1.2	Dynamic Programming - Theory . . . . .	5
2.1.3	Choice based linear programming . . . . .	6
2.2	Approximate Dynamic Programming . . . . .	7
2.2.1	Approximate Policy Iteration . . . . .	7
2.2.2	Greedy Heuristic to determine offerset . . . . .	8
2.2.3	Update of parameters . . . . .	9
2.2.4	More stuff . . . . .	10
<b>3</b>	<b>Examples</b>	<b>11</b>
3.1	Example0 . . . . .	11
<b>4</b>	<b>Main chapter</b>	<b>13</b>
4.1	Notation . . . . .	13
4.1.1	Examples . . . . .	13
4.2	Implementation . . . . .	13
4.2.1	Network flight example . . . . .	14
<b>5</b>	<b>Comparison of different policies</b>	<b>15</b>
5.1	Two sample test . . . . .	15
<b>6</b>	<b>Thoughts on implementation</b>	<b>16</b>
<b>7</b>	<b>Interesting findings</b>	<b>17</b>
7.1	DP offer nothing at start . . . . .	17
	<b>Bibliography</b>	<b>17</b>
	<b>Bibliography</b>	<b>20</b>
	<b>Nomenclature</b>	<b>21</b>

# Chapter 1

## Aufzuschreiben

This is a quick overview of my results, that I want to incorporate in the thesis.

1. Computing the Offer Set Dynamically as presented in Bront et Al 5.2 and reproduce example of Table 3 to proof correct implementation C:

Users

Stefan

LRZ Sync+Share

Masterarbeit-Klein

Code

notebooks

0427 - Deep dive into DPD Offer Set Calculation.ipynb

2. Most structure of all coding until May also in C:

Users

Stefan

LRZ Sync+Share

Masterarbeit-Klein

Code

notebooks

doc\_Koch.ipynb

# Chapter 2

## Problem and Approaches to Solution

### 2.1 The problem

Here, we want to lay out the classical revenue management problem and the “brute force” approach to solve.

Consider a firm that produces products  $j = 1, \dots, n$  with revenues  $\mathbf{r} = (r_1, \dots, r_n)^T$ . Resources  $h = 1, \dots, m$  are used for production. In order to produce one unit of product  $j$ , resources  $\mathbf{a}_j = (a_{1j}, \dots, a_{mj})^T$  are necessary, with  $a_{hj} = 1$  if resource  $h$  is needed for production of product  $j$  and  $a_{hj} = 0$  otherwise. Initially, the capacity is described by  $\mathbf{c}^0 = (c_1^0, \dots, c_m^0)^T$ .

The booking horizon is modelled by sufficiently small time periods  $t = 0, \dots, T - 1$ , such that in each time period at most one customer arrives. This customer also purchases at most one product. If product  $j$  is purchased at time  $t$ , the capacity reduces to  $\mathbf{c}^{t+1} = \mathbf{c}^t - \mathbf{a}_j$ . Time moves forward, such that the last selling might occur at time  $T - 1$ .

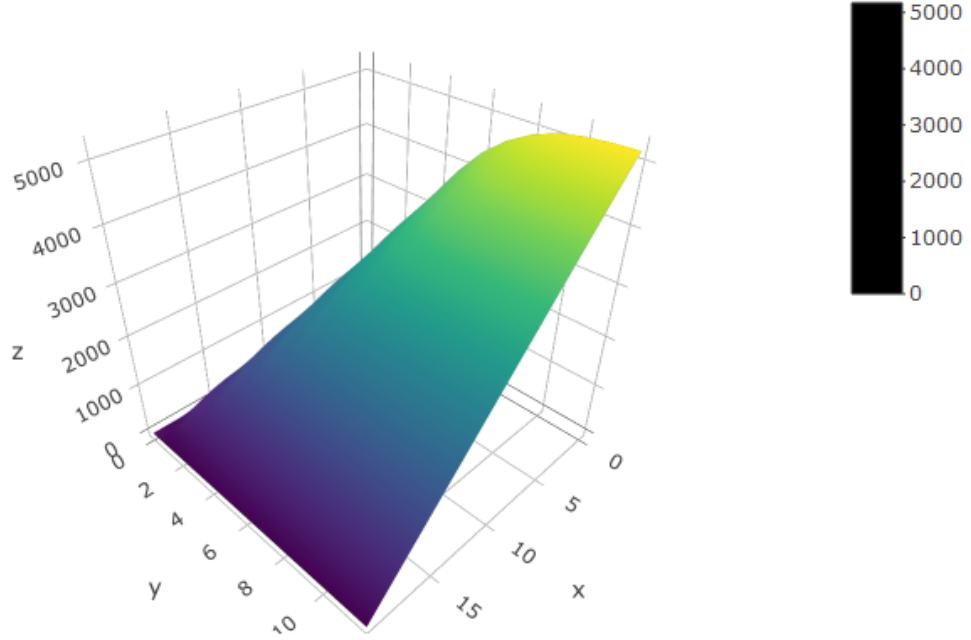
The firm wants to increase the value of the products sold and has flexibility in the sets offered. Thus, the decision variables at each time point  $t$  are given by  $\mathbf{x}^t = (x_1^t, \dots, x_n^t)^T$  with  $x_j^t = 1$  if product  $j$  is offered at time  $t$  and  $x_j^t = 0$  otherwise. Thus, at each time point  $t$  and each capacity  $\mathbf{c}$ , the offer set  $\mathbf{x}$  has to be determined. Note that in machine learning terms, this mapping is defined as a policy  $\boldsymbol{\psi} = (\psi^1, \dots, \psi^T)^T$ .

One popular method of describing the probabilities of purchases is to have each customer belonging to one customer segment  $l = 1, \dots, L$ , each of which following a multinomial logit model (MNL). A customer of segment  $l$  arrives with probability  $\lambda_l$ . His preference weights are given by  $\mathbf{u}_l = (u_{l1}, \dots, u_{ln})^T$  and no purchase preference of  $u_{l0}$ . Note:  $u_{lj} > 0$  if consumer of segment  $l$  might purchase product  $j$  (the higher, the more interested) and  $u_{lj} = 0$  if customer is not interested in product. The probability of purchasing product  $j$  when set  $\mathbf{x}$  is offered is given by  $p_{lj}(\mathbf{x}) = \frac{u_{lj}x_j}{u_{l0} + \sum_{p \in [n]} u_{lp}x_p}$  and the no-purchase probability is given by  $p_{l0}(\mathbf{x}) = 1 - \sum_{j \in [n]} p_{lj}$ . Together with the uncertainty of which customer segment arrives (if any), we end up at a purchase probability for product  $j$  given  $\mathbf{x}$  of  $p_j(\mathbf{x}) = \sum_{l \in [L]} \lambda_l p_{lj}(\mathbf{x})$  and a no purchase probability of  $p_0(\mathbf{x}) = 1 - \sum_{j \in [n]} p_j(\mathbf{x})$ . Thus,  $p_j(\mathbf{x})$  can be interpreted as the deterministic quantity of product  $j$  being sold, if set  $\mathbf{x}$  is offered.

#### 2.1.1 Dynamic Programming

The first solution approach might be to maximize the expected value of all revenues to gain given period  $t$  and capacity  $\mathbf{c}$ , denoted by the value function  $V^t(\mathbf{c})$ . This can then be computed recursively as

Figure 2.1: Value function of single leg flight example with  $t \in [20]$  on x-axis,  $c \in [12]$  on y-axis.



$$V^t(\mathbf{c}) = \max_{x^t \in \{0,1\}^n} \left\{ \sum_{j \in [n]} p_j(x_t) (r_j + V^{t+1}(\mathbf{c} - \mathbf{a}_j)) + p_0 V^{t+1}(\mathbf{c}) \right\} \quad (2.1) \quad \{\text{eq-Bellman}\}$$

$$= \max_{x^t \in \{0,1\}^n} \left\{ \sum_{j \in [n]} p_j(x_t) (r_j - \Delta_j V^{t+1}(\mathbf{c})) \right\} + V^{t+1}(\mathbf{c}) \quad \forall t, \mathbf{c} \geq 0 \quad (2.2)$$

with  $\Delta_j V^{t+1}(\mathbf{c}) := V^{t+1}(\mathbf{c}) - V^{t+1}(\mathbf{c} - \mathbf{a}_j)$  and boundary conditions  $V^{T+1}(\mathbf{c}) = 0$  if  $\mathbf{c} \geq \mathbf{0}$  and  $V^t(\mathbf{c}) = -\infty$  if  $\mathbf{c} \not\geq \mathbf{0}$ .

The optimal offerset  $\mathbf{x}^t$  is then given by

$$\mathbf{x}^t = \arg \max_{x^t \in \{0,1\}^n} \left\{ \sum_{j \in [n]} p_j(x_t) (r_j - \Delta_j V^{t+1}(\mathbf{c})) \right\} + V^{t+1}(\mathbf{c}) \quad \forall t, \mathbf{c} \geq 0 \quad (2.3) \quad \{\text{eq-DP-opt0}\}$$

The resulting value function can be plotted in a three dimensional plot as done in Figure 2.1.

### 2.1.2 Dynamic Programming - Theory

As this whole thesis is based on the ideas of Dynamic Programming, we want to give a short overview of the underlying mathematical theory and directly combine it to our setting.

Dynamic Programming (DP) refers to a broad collection of algorithms used to compute optimal policies given a perfect model of the environment as a Markov Decision Process (MDP), as stated in Sutton and Barto (2018). A MDP is characterized by a state set  $\mathcal{S}$  (the varying capacities  $\mathbf{c}$  together with the current point in time  $t$ ), an action set  $\mathcal{A}$  (the offer sets  $\mathbf{x}$ ) and reward sets  $\mathcal{R}$  (the revenue  $r$  if a product is sold). A MDP evolves over time, but for the evolution from time  $t$  to  $t + 1$  only the current state  $s_t$  is relevant and the history of previous states  $s_h, h < t$  can be ignored. The dynamic is given by a set of probabilities  $p(s', r | s, a)$  (capacities reduce according to which random product the random customer has bought according to his preferences).

The goal of DP is to determine the optimal policy, i.e. which action to choose at each given state. Key to the solution is the usage of value functions as seen above. These fulfil the Bellman optimality equations as stated in Equation (2.1).

In our setting, uniqueness of the optimal value function is ensured, as the MDP is guaranteed to terminate under any policy because time is moving forward no matter which action we choose, compare Sutton and Barto (2018). Thus, while the problem can be solved, due to large scale decision problems and the curse of dimensionality, we are lacking computing power to solve it exactly and have to come up with approximate solution methods.

One approximate solution method is the usage of Approximate Policy Iteration (API), which consists of two steps. The policy evaluation step (inner loop Line 4 to Line 15) evaluates a fixed policy over a set of sample paths  $\omega_i, i = 1, \dots, I$ . The policy improvement step (outer loop Line 2, Line 3, Line 16) improves the policy over the iterations  $k = 1, \dots, K$ .

### 2.1.3 Choice based linear programming

based on Bront et al and van Ryzin.

Taking an eagle-eye perspective, the company has to decide which sets to offer in every single time period, but as the probabilities don't change over time, the specific time point when to offer one set is indistinguishable, we can aggregate over time. Let us introduce a few more notation to get a hand on this. Let  $R(\mathbf{x})$  represent the expected revenue given offer set  $\mathbf{x}$ , i.e.

$$R(\mathbf{x}) = \sum_{j \in [n]} r_j p_j(\mathbf{x}) \quad . \quad (2.4)$$

Furthermore, let  $\mathbf{Q}(\mathbf{x}) = (Q_1(\mathbf{x}), \dots, Q_m(\mathbf{x}))^T$  denote the expected capacity being used, i.e.

$$Q_h(\mathbf{x}) = \sum_{j \in [n]} a_{hj} p_j(\mathbf{x}) \quad . \quad (2.5)$$

We can denote the total time set  $\mathbf{x}$  is offered by  $t(\mathbf{x})$  and the set of all possible offersets as  $N = \{0, 1\}^n$ . Thus, we can formulate the choice-based linear program (CDLP) as presented in van Ryzin and Liu. Here the optimization problem is solved by varying the decision variables  $t(\mathbf{x})$ .

$$V^{CDLP} = \max \sum_{\mathbf{x} \in N} R(\mathbf{x})t(\mathbf{x}) \quad (2.6)$$

$$\text{s.t. } \sum_{\mathbf{x} \in N} Q_h(\mathbf{x})t(\mathbf{x}) \leq c_h, \forall h \quad (2.7)$$

$$\sum_{\mathbf{x} \in N} t(\mathbf{x}) \leq T, \quad (2.8)$$

$$t(\mathbf{x}) \geq 0 \quad \forall \mathbf{x} \in N. \quad (2.9)$$

Korrektheit getestet mit C:\Users\Stefan\LRZ\_Sync+Share\Masterarbeit-Klein\Code\H\_CDLP\_single-Example0 in settings Ergebnisse in C:\Users\Stefan\LRZ\_Sync+Share\Masterarbeit-Klein\Code\Resu.

## 2.2 Approximate Dynamic Programming

### 2.2.1 Approximate Policy Iteration

ADP builds upon Equation (2.3) and approximates opportunity costs additively using bid prices  $\pi_h(t, c_h)$ , i.e.  $\Delta_j V^{t+1}(\mathbf{c}) = \sum_{h \in [m]} a_{hj} \pi_h(t+1, c_h)$ . Thus, the policy is given as the solution of

$$\pi \mathbf{x}^t = \arg \max_{\mathbf{x}^t \in \{0,1\}^n} \left\{ \sum_{j \in [n]} p_j(x_t) \left( r_j - \sum_{h \in [m]} a_{hj} \pi_h(t+1, c_h) \right) \right\} \quad (2.10)$$

Thus, the policy is completely described by the bid prices  $\pi_h(t, c_h) \forall t$ .

Let us discuss the algorithm in depth. In the beginning (Line 1), all parameters are set to zero which results in the greedy policy of offering the set resulting in the highest expected revenue (no consideration of costs).

In each policy iteration  $k$ , the current policy is evaluated first. To do so, we use the current policy to evaluate  $I$  random paths. Line 3 sets up the dataset for revenues to go  $\hat{V}_t^i$  and remaining capacities  $\hat{\mathbf{C}}_t^i$  for all times  $t$  and sample paths  $i$ .

For each sample path  $i$ , a storage unit for revenue  $\hat{r}_t$  and capacities  $\hat{\mathbf{c}}$  is created for all times  $t$  (Line 5). Furthermore, the capacity is set to the initial capacity in Line 6.

For each point in time, the current capacity  $\hat{\mathbf{c}}_t$  is stored in Line 8, the current bid prices are calculated in Line 9 and used to determine the offerset  $\mathbf{x}$  in Line 10. More details on the determination of the offerset can be found in Section 2.2.2. With, this, a sales event  $j'$  is simulated and in case a product has been sold, the revenue is stored and capacity adjusted (note: the revenue thus belongs to  $t$  and capacity will affect  $t+1$ ) (Line 11 to Line 13).

After simulating one sample path, the value function for this sample path  $\hat{V}_t^i$  is updated for all times  $t$  in Line 14, and so are the capacities in line Line 15.

With all the information from the  $I$  sample paths the parameters  $(\theta_t, \pi_t)$  are updated. Note that the old parameters are used as starting values and the current number of policy iterations  $k$  is also passed to potentially take care of exponential smoothing. Thus, we have  $(1+h) * T$  parameters ( $\theta_t$  and  $\pi_t$ ) and  $2TI$  data points. More details on the update of the parameters can be found in Section 2.2.3.

Actions possible in  $t = 0, \dots, T-1$ .

Overview of parameters:

**Algorithm 1** Approximate policy iteration

1: Set $\theta_t = 0$ and $\pi_t = \mathbf{0} \forall t = 0, \dots, T - 1$	{alg-API}
2: <b>for</b> $k = 1$ <b>to</b> $K$ <b>do</b>	{alg-API1}
3:     Set $\hat{V}_t^i = 0$ and $\hat{\mathbf{C}}_t^i = 0 \forall t = 0, \dots, T - 1, \forall i = 1, \dots, I$	{alg-API-Pit}
4: <b>for</b> $i = 1$ <b>to</b> $I$ <b>do</b>	{alg-API-Pit}
5:         Set $\hat{r}_t = 0$ and $\hat{\mathbf{c}}_t = 0 \forall t = 0, \dots, T - 1$	{alg-API-Pev}
6:         Initialize $\mathbf{c} = \mathbf{c}^0$	{alg-API5}
7: <b>for</b> $t = 0$ <b>to</b> $T-1$ <b>do</b>	{alg-API6}
8: $\hat{\mathbf{c}}_t := \mathbf{c}$	{alg-API8}
9:             Compute $\pi(t, \mathbf{c})$	{alg-API9}
10:            Compute $\mathbf{x} = \text{determineOfferset}(\pi(t, \mathbf{c}), \epsilon_t)$	{alg-API9}
11:            Simulate a sales event $j' \in \{0, 1, \dots, n\}$	{alg-API10}
12: <b>if</b> $j' \in \{1, \dots, n\}$ <b>then</b>	{alg-API11}
13: $\hat{r}_t = r_{j'}$ and $\mathbf{c} = \mathbf{c} - \mathbf{a}_{j'}$	{alg-API13}
14:            Compute $\hat{V}_t^i = \sum_{\tau=t}^T \hat{r}_\tau \quad \forall t = 0, \dots, T - 1$	{alg-API14}
15:            Assign $\hat{\mathbf{C}}_t^i = \hat{\mathbf{c}}_t \quad \forall t = 0, \dots, T - 1$	{alg-API14}
16: $(\theta_t, \pi_t) = \text{updateParameters}(\hat{V}_t^i, \hat{\mathbf{C}}_t^i, \theta_t, \pi_t, k) \quad \forall t = 0, \dots, T - 1, \forall i = 1, \dots, I$	{alg-API14}
<b>return</b> $(\theta_t, \pi_t) \quad \forall t = 1, \dots, T$	{alg-API-Bpd}

1.  $\theta_t$  optimization parameter (offset)
2.  $\pi_t$  optimization parameter (bid price for each resource)
3.  $\hat{V}_t = 0$  all sample revenues to go for each sample for each time
4.  $\hat{\mathbf{C}}_t = 0$  all sample available capacities for each sample for each resource for each time
5.  $r_t$  sample revenue generated at time  $t$
6.  $\mathbf{c}$  available capacities for each resource at current time
7.  $\mathbf{c}^0$  starting capacities
8.  $\mathbf{x}$  offerset at current time
9.  $\epsilon_t$  epsilon used at current time

### 2.2.2 Greedy Heuristic to determine offerset

The following algorithm is based on the ideas of the greedy heuristic for the column generation subproblem outlined in Bront et al. (2009).

Our goal is to determine a reasonable set of products to offer in a fast manner. Thus, we use a heuristic and cut down the amount of products to consider as fast as possible.

The function  $\text{determineOfferset}(\pi, \epsilon)$  calculates the set to offer depending on the current bid prices  $\pi$  via the greedy algorithm layed out in Bront et al. (2009). To account for the exploration vs exploitation dilemma, an epsilon-greedy strategy is used. With a probability of  $\epsilon/2$  either no product is offered at all or all products with positive contribution  $r_j - \sum_{h \in [m]} a_{hj} \cdot \pi_h$  are offered. With a probability of  $1 - \epsilon$ , the proper calculated set is offered.



Value Funktion gebündelt dargestellt und ueberall mit  $\lambda$ . Vergleiche zu Bront et al. 4.2.2, wo in 3. ohne  $\lambda$  und in 4.a mit  $\lambda$ .

---

**Algorithm 2** Greedy Heuristic

---

```

1: Value( $X$ ) :=  $\sum_{l=1}^L \lambda_l \frac{\sum_{i \in X} (r_i - A_i^T \pi) u_{li}}{\sum_{i \in X} u_{li} + u_{l0}}$ 
2:  $S := \emptyset$ ,  $S' := \{j \in N : r_j - A_j^T \pi > 0\}$ 
3:  $j^* := \arg \max_{j \in S'} \text{Value}(\{j\})$ 
4: repeat
5:    $S := S \cup \{j^*\}$ ,  $S' := S' \setminus \{j^*\}$ 
6:    $j^* := \arg \max_{j \in S'} \text{Value}(S \cup \{j\})$ 
7: until  $\text{Value}(S \cup \{j^*\}) \leq \text{Value}(S)$ 
8: return  $S$ 

```

---

{alg-GreedyH}

{alg-L1}

Let  $S'$  be the set of products with positive reduced costs, i. Line 2

### 2.2.3 Update of parameters

{sec-updateP}

The function  $(\theta_t, \pi_t) = \text{updateParameters}(\hat{V}_t, \hat{\mathbf{C}}_t, \theta_t, \pi_t, k)$  really optimizes the following least squares optimization problem for all parameters  $(t = 1, \dots, T)$  at the same time.

$$V_t(\theta_t, \pi_t, \mathbf{c}_t) := \theta_t + \sum_{h=1}^m \sum_{s=1}^{S_h} \pi_{ths} f_{hs}(c_h) \quad (2.11)$$

$$f_{hs}(c_h) := \begin{cases} 0 & \text{if } c_h \leq b_h^{s-1} \\ c_h - b_h^{s-1} & \text{if } b_h^{s-1} < c_h \leq b_h^s \\ b_h^s - b_h^{s-1} & \text{if } b_h^s < c_h \end{cases} \quad (2.12) \quad \{\{\text{def-f}\}\}$$

Equation (2.12) describes the occupied amount of capacity of interval  $(b_h^{s-1}, b_h^s]$ .

The following optimization problem depends on the old parameters  $\theta_t = \theta_t^k$  and  $\pi_t = \pi_t^k$  to determine the optimal parameter  $\theta_t^{\text{update}}$  and  $\pi_t^{\text{update}}$ .

$$\min \sum_{i=1}^I \sum_{t=1}^T \left( \hat{V}_t^i - V_t(\theta_t, \pi_t, \mathbf{c}_t^i) \right)^2 \quad (2.13)$$

$$s.t. \quad (2.14)$$

$$\theta_t \geq 0 \quad \forall t \quad (2.15)$$

$$\max_{j=1, \dots, n} r_j \geq \pi_{ths} \geq 0 \quad \forall t, h, s \quad (2.16)$$

$$\pi_{ths} \geq \pi_{th,s+1} \quad \forall t, h, s = 1, \dots, S_h - 1 \quad (2.17)$$

$$\theta_t \geq \theta_{t+1} \quad \forall t = 1, \dots, T - 1 \quad (2.18)$$

$$\pi_{ths} \geq \pi_{t+1,hs} \quad \forall t = 1, \dots, T - 1 \quad (2.19)$$

The final parameters  $\theta_t^{K+1}$  and  $\pi_t^{K+1}$  can obtained via two possible equally possible ways. One is the so called exponential smoothing, where in each iteration  $k$  the parameter for the

next iteration  $k + 1$  is calculated via:

$$\theta_t^{k+1} = \left(1 - \frac{1}{k}\right) \theta_t^k + \frac{1}{k} \theta_t^{\text{update}} \quad (2.20)$$

$$\pi_t^{k+1} = \left(1 - \frac{1}{k}\right) \pi_t^k + \frac{1}{k} \pi_t^{\text{update}} \quad (2.21)$$

The other one uses  $\theta_t^{k+1} = \theta_t^{\text{update}}$  and  $\pi_t^{k+1} = \pi_t^{\text{update}}$  and averages at the very end.

$$\theta_t^{K+1} = \frac{1}{K} \sum_{k=1}^K \theta_t^k \quad (2.22)$$

$$\pi_t^{K+1} = \frac{1}{K} \sum_{k=1}^K \pi_t^k \quad (2.23)$$

**Proof.**

Den Beweis sauber ausfuehren. Hier sind die Inhalte. Die Aussage stimmt, falls die gefundenen optimalen Loesungen in jeder Iteration stets dieselben sind. Dies ist der Fall, wenn es stets nur ein Minimum gibt. Wir haben eine quadratische Zielfunktion, die sozusagen eine mehrdimensionale, nach oben geoeffnete Parabel zeigt, die genau ein globales Minimum besitzt. Weitere noetige Punkte: eindeutiges globales Minimum ueber positiv definite Hesse-Matrix. Optimum auch zulaessig (schwierig?)

■

## 2.2.4 More stuff

Calculations:  $\pi(t, \mathbf{c}) = \pi_h(t, c_h)$  for  $h \in [m]$

$$\pi_h(t, c_h) = \begin{cases} \infty & \text{if } c_h = 0 \\ \sum_{s=1}^{S_h} \pi_{ths} \mathbb{1}_{(b_h^{s-1}, b_h^s]}(c_h) & \text{otherwise.} \end{cases} \quad (2.24)$$

$$(2.25)$$

Fuer Line 9 verwende Zeit  $\mathbf{t}$  statt  $\mathbf{t}+1$ . Grund: Kenne Informationen zur Zukunft nicht.

A sales event is simulated by first having one or zero customer arrive at random. In case a customer arrives, its preference function given the offer set determines the probability according to which one product is sold ( $j' \in \{1, \dots, n\}$ ) or no product is sold ( $j' = 0$ ).

# Chapter 3

## Examples

### 3.1 Example0

One interesting result is that for exact reproduction of the results of Bront et al for their example0, the empty offerset has to be excluded from the optimization problem. Otherwise (including empty set) results in the same optimal value but via another combination of the decision variables.

*Example 0* represents a small airline network. Three cities are interconnected by three flights (legs), with a capacity vector  $\mathbf{c} = (10, 5, 5)^T$ . The booking horizon consists of  $T = 30$  periods and there are five customer segments with preferences as in Table 3.1

Table 3.1: Segment definition for Example 0<sup>tb-Example0-Customers</sup>

Segment	$\lambda_l$	Consideration set	Preference vector	No purchase preference	Description
1	0.15	$\{1, 5\}$	$(5, 8)$	2	Price sensitive, nonstop

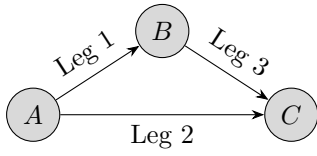


Figure 3.1: Airline network for example 0. <sup>fig-Example0</sup>

Table 3.2: Airline network for example 0 (products). <sup>tb-Example0-Products</sup>

Product	Origin-destination	Fare
1	$A \rightarrow C$	1,200
2	$A \rightarrow B \rightarrow C$	800
3	$A \rightarrow B$	500
4	$B \rightarrow C$	500
5	$A \rightarrow C$	800
6	$A \rightarrow B \rightarrow C$	500
7	$A \rightarrow B$	300
8	$B \rightarrow C$	300

# Chapter 4

## Main chapter

### 4.1 Notation

1. time index always at the top
2. index for products always at bottom

#### 4.1.1 Examples

##### Single-leg flight example

For reasons of comparability, we use the same example as in Koch (2017). An airline offers four products with revenues  $\mathbf{r} = (1000, 800, 600, 400)^T$  over  $T = 400$  periods. Only one customer segment exists with arrival probability of  $\lambda = 0.5$  and preference weights  $\mathbf{u} = (0.4, 0.8, 1.2, 1.6)^T$ . Different network loads can be analyzed by varying initial capacity  $c^0 \in \{40, 60, \dots, 120\}$  and varying no-purchase preference weights  $u_0 \in \{1, 2, 3\}$ .

### 4.2 Implementation

Here, we present the results of the exact calculation for the single leg flight example.

Storage folder:

"C:/Users/Stefan/LRZ Sync+Share/Masterarbeit-Klein/Code/Results/singleLegFlight-True-DP-

Log:

Time (starting): 2019-06-11 09:17:42.448064

Example : singleLegFlight

use\_var\_capacities : True

Total time needed:

151.31108283996582 seconds =

2.521851380666097 minutes

Results:

	capacity	no-purchase preference	DP-value	DP-optimal offer set
0	40	1	35952	(1, 1, 0, 0)
1	60	1	49977.1	(1, 1, 0, 0)
2	80	1	59969.8	(1, 1, 1, 0)
3	100	1	65956.4	(1, 1, 1, 0)
4	120	1	68217.2	(1, 1, 1, 1)
5	40	2	35952	(1, 1, 0, 0)
6	60	2	49977.1	(1, 1, 0, 0)
7	80	2	59969.8	(1, 1, 1, 0)
8	100	2	65956.4	(1, 1, 1, 0)
9	120	2	68217.2	(1, 1, 1, 1)
10	40	3	35952	(1, 1, 0, 0)
11	60	3	49977.1	(1, 1, 0, 0)
12	80	3	59969.8	(1, 1, 1, 0)
13	100	3	65956.4	(1, 1, 1, 0)
14	120	3	68217.2	(1, 1, 1, 1)

Furthermore, we want to visualize the value function and its approximations. We use the example of  $c^0 = 60$  and  $u_0 = 1$ .

#### 4.2.1 Network flight example

# Chapter 5

## Comparison of different policies

### 5.1 Two sample test

Our goal is to evaluate two different policies. One policy is considered better than another policy if it generally leads to higher total revenues. To put this vague phrasing into a scientifically sound comparison, we use an approximate two sample test, which we'll describe here in more detail. We base this overview and notation mainly on Ch. 14.7 of Bamberg et al. (2011) with extensions as in Ch. 11.3 of Fahrmeir et al. (2007).

We have a total of  $M$  sample paths each consisting of  $T$  time steps, i.e.  $M \times T$  random numbers determining which customer arrives and  $M \times T$  random numbers leading to which product is purchased. Those two numbers are called *exogenous information*. When applying policy  $A$ , for each sample path at each time step, the respective offerset is presented to the current customer and its preferences determine which product is purchased. All purchases of one sample path are aggregated and stored, such that policy  $A$  results in the vector of values  $\mathbf{v}^A \in \mathbb{R}^M$ . Accordingly,  $\mathbf{v}^B$  is determined. As both policies rely on the same exogenous information, the samples are dependent and a paired samples t-test has to be used.

Thus, we end up with two observations of size  $M$

$$v_1^A, \dots, v_M^A \text{ respectively } v_1^B, \dots, v_M^B .$$

Let  $\bar{\mathbf{v}}^A$  and  $S_A^2$  resp.  $\bar{\mathbf{v}}^B$  and  $S_B^2$  be the empirical mean and empirical variance. Furthermore,  $\mu^A$  and  $\sigma_A^2$  denote the expected value and variance of  $v_i^A$ , respectively  $\mu^B$  and  $\sigma_B^2$  denote the expected value and variance of  $v_i^B$ . The hypothesis that policy  $A$  is better than policy  $B$  results in the following null and alternative hypothesis:

$$H_0 : \mu^A \leq \mu^B, H_1 : \mu^A > \mu^B$$

As  $\bar{\mathbf{v}}^A$  and  $\bar{\mathbf{v}}^B$  are unbiased estimators of  $\mu^A$  and  $\mu^B$ , the difference  $D = \bar{\mathbf{v}}^A - \bar{\mathbf{v}}^B$  can be used to test the hypothesis. A large value of  $D$  represents the data to be in favour of  $H_1$ .

As we don't know the distribution of  $v_i^A$  or  $v_i^B$  and the sample size  $M > 30$ , we know for the test statistic

$$T = \frac{\bar{\mathbf{v}}^A - \bar{\mathbf{v}}^B}{\sqrt{\frac{S_A^2 + S_B^2}{M}}} \sim N(0; 1) .$$

With this knowledge, we can compute the  $p$ -value according to  $p = 1 - \Phi(T)$  with  $\Phi(\cdot)$  being the standard normal distribution function. Note: The  $p$ -value can be seen as evidence against  $H_0$ . A small  $p$ -value represents strong evidence against  $H_0$  and the null hypothesis should be rejected if the  $p$ -value is below a significance threshold  $\alpha$ .

## Chapter 6

### Thoughts on implementation

This shall be a short overview of my thoughts regarding the implementation:

1. I wanted to produce reproducible results, i.e. my code shall be usable on a larger scope. Logic and Settings shall be separated. Completed code shall be script based and after run, create a folder with its running configuration.
2. Be careful, when storing list or dataframes, as a deepcopy might have to be used.
- 3.



# Chapter 7

## Interesting findings

### 7.1 DP offer nothing at start

In the single leg flight scenario with no purchase preference of 1, we found that it makes no difference (due to rounding error), whether just the most expensive product is offered or no product at all. Thus, we want to quickly note this result down here.

We want to estimate the actions in time  $t = 1$ . After having the final results computed for all capacities in  $t = 2$ , we arrive at the following values:

Furthermore, we have the following offersets:

resulting in the purchase probabilities:

This results in the following expected values,

[1000.815.38461538890.90909091810.52631579955.55555556835.29411765893.33333333826.086956521000.84

i.e. offering nothing and offering just the most expensive product both result in an expected value of 1000.

As it almost never makes sense to offer no product, we switched the offer set with no products to offer at the very end. The python function `numpy.argmax` returns the index of the first maximum (if there are more than one).

	value	offer_set_optimal	num_offer_set_optimal
0	0.0	0	0
1	1000.0	0	2
2	2000.0	0	2
3	3000.0	0	2
4	4000.0	0	2
5	5000.0	0	2
6	6000.0	0	2
7	7000.0	0	2
8	8000.0	0	2
9	9000.0	0	2
10	10000.0	0	2
11	11000.0	0	2
12	12000.0	0	2
13	13000.0	8	1
14	14000.0	8	1
15	15000.0	8	1
16	16000.0	8	1
17	17000.0	8	1
18	18000.0	8	1
19	19000.0	8	1
20	20000.0	8	1

	0	1	2	3
0	0	0	0	0
1	0	0	0	1
2	0	0	1	0
3	0	0	1	1
4	0	1	0	0
5	0	1	0	1
6	0	1	1	0
7	0	1	1	1
8	1	0	0	0
9	1	0	0	1
10	1	0	1	0
11	1	0	1	1
12	1	1	0	0
13	1	1	0	1
14	1	1	1	0
15	1	1	1	1

	0	1	2	3	4
0	0.000000	0.000000	0.000000	0.000000	1.000000
1	0.000000	0.000000	0.000000	0.307692	0.692308
2	0.000000	0.000000	0.272727	0.000000	0.727273
3	0.000000	0.000000	0.157895	0.210526	0.631579
4	0.000000	0.222222	0.000000	0.000000	0.777778
5	0.000000	0.117647	0.000000	0.235294	0.647059
6	0.000000	0.133333	0.200000	0.000000	0.666667
7	0.000000	0.086957	0.130435	0.173913	0.608696
8	0.142857	0.000000	0.000000	0.000000	0.857143
9	0.066667	0.000000	0.000000	0.266667	0.666667
10	0.076923	0.000000	0.230769	0.000000	0.692308
11	0.047619	0.000000	0.142857	0.190476	0.619048
12	0.090909	0.181818	0.000000	0.000000	0.727273
13	0.052632	0.105263	0.000000	0.210526	0.631579
14	0.058824	0.117647	0.176471	0.000000	0.647059
15	0.040000	0.080000	0.120000	0.160000	0.600000

# Bibliography

- Bamberg, G., Baur, F., and Krapp, M. (2011). *Statistik*. Lehr- und Handbücher der Wirtschafts- und Sozialwissenschaften. Oldenbourg, München, 16., überarb. Aufl. edition.
- Bront, J. J. M., Méndez-Díaz, I., and Vulcano, G. (2009). A column generation algorithm for choice-based network revenue management. *Operations Research*, 57(3):769–784.
- Fahrmeir, L., Künstler, R., Pigeot, I., and Tutz, G. (2007). *Statistik: Der Weg zur Datenanalyse*. Springer-Lehrbuch. Springer, Berlin, 6., überarb. Aufl. edition.
- Koch, S. (2017). Least squares approximate policy iteration for learning bid prices in choice-based revenue management. *Computers & Operations Research*, 77:240–253.
- Sutton, R. S. and Barto, A. (2018). *Reinforcement learning: An introduction*. Adaptive computation and machine learning. The MIT Press, Cambridge, MA and London, second edition edition.

# Nomenclature

**API** approximate policy iteration

**CDLP** choice-based linear program

**DP** dynamic programming

**MDP** markov decision process

**MNL** multinomial logit model