

Research on Production Decision-Making Based on Dynamic Programming

Leyi Huang *

School of information science and technology, Jinan University, Guangzhou, China, 511443

* Corresponding Author Email: huangleiyi450@gmail.com

Abstract. For the current enterprise in the production process faces the problem of high scrap rate, high inspection cost and low profit, how to optimize the decision-making in the multi-stage production process has become an urgent issue. This study focuses on innovating the decision-making optimization model of the multi-stage production process, and firstly, the production process of the electronics factory is clearly defined into four stages: part inspection, finished product inspection, defective product disassembly and disassembly of parts after disassembly. In this paper, we innovatively propose to combine the interdependence and feedback mechanism between different stages in the production process, use dynamic programming and economic benefit evaluation models to optimize the cost of each stage, and further construct the Markov decision process by combining the reward function and state transition probability, and use the Bellman equation to calculate the optimal decision scheme and verify its accuracy. This study provides an effective new fusion model for decision-making optimization in the production process, and makes a breakthrough in multi-stage joint optimization, and finds that strict quality control, especially in the semi-finished product stage, can significantly reduce costs, which provides guidance for enterprises to improve economic benefits. To reduce production costs and improve corporate profits.

Keywords: Production Decisions; Dynamic Programming; Economic Efficiency Assessment Model; Markov Decision Process.

1. Introduction

As the demand for electronic products continues to grow, manufacturing companies are under increasing pressure to control quality and manage costs. During the production process, any failure of spare parts may lead to a substandard final product, and the uncertainty in the assembly process also puts the final product at risk of defects. Therefore, how to optimize inspection and disassembly decisions to achieve quality assurance while reducing costs has become a key challenge in production management.

Zhang R et al. [1] proposed a dynamic stochastic production decision model based on analytic-simulation feedback. Firstly, the static stochastic analytic model is solved for the optimal decision and input into the simulation model. Then, the relaxed decision parameters are brought into the analytic model to model and solve the problem again in a loop and simulate until it meets the stopping criterion. Yang C et al. [2] proposed dynamic planning for multi-stage assembly line operation scheduling, algorithm design, and analysis. Cai M et al. [3] used dynamic programming to solve the optimal production strategy problem with Lingo implementation.

Dynamic programming [4] has great advantages in dealing with the optimization problem of staged decision-making, and the core "optimal principle" can balance the cost and benefit of the production process, flexibly plan the production strategy, and has good scalability, so this model is followed in this paper. However, the existing research fails to fully consider the interdependence and feedback mechanism between different stages of the production process, and lacks the comprehensive optimization of each link in the complex production process. This paper proposes a multi-stage production optimization framework combining dynamic programming and economic benefit evaluation model, which further combines the Markov decision-making process and innovates the balance between quality control and cost management. It not only solves the complexity of multi-

stage production decision-making, but also has higher flexibility and adaptability of the model, which promotes the rationalization of production resource allocation.

2. Dynamic planning economic efficiency assessment modeling

2.1. Problem assumption

The specific processes in the production process of the present enterprise are as follows:

- 1) Whether the parts (Part 1 and/or Part 2) are tested or not, if a part is not tested, it goes directly to the assembly; otherwise, the detected nonconforming part is discarded;
- 2) Whether to test each piece of assembled finished product, if not, the assembled finished product directly into the market; otherwise only qualified finished product into the market:
- 3) Whether the detected unqualified finished products are disassembled, if not, directly discard the unqualified finished products, otherwise repeat steps (1) and (2) for the disassembled spare parts;
- 4) For non-conforming products purchased by users, the enterprise will exchange them unconditionally and incur certain exchange losses (such as logistics costs, enterprise reputation, etc.). Repeat step (3) for returned nonconforming products.

The data for cases of two known spare parts and finished products are shown in Table 1 and 2 below:

Table 1. Production Conditions1

Situation	Component1		Component2		product		
	defect rate	Inspection cost (\$)	defect rate	Inspection cost (\$)	defect rate	Assembly cost (\$)	Inspection cost (\$)
1	10%	2	10%	3	10%	6	3
2	20%	2	20%	3	20%	6	3
3	10%	2	10%	3	10%	6	3
4	20%	1	20%	1	20%	6	2
5	10%	8	20%	1	10%	6	2
6	5%	2	5%	3	5%	6	3

Table 2. Defective Product

Situation	1	2	3	4	5	6
Replacement loss	6	6	30	30	10	10
Disassembly cost	5	5	5	5	5	40

2.2. Model preparation

0-1 planning [5] is a special kind of integer planning problem in which the decision variables can only take 0 or 1, which is suitable for problems that require binary decision making. Therefore, 0-1 planning is used to represent the state of each stage when making a decision on whether to proceed with each stage of the production process of a company. The production stages of the electronics factory are now divided into four:

- 1) Stage 1: Part 1 and Part 2 detection decision; state representation: (S_1, S_2)
 $S_1 \in \{0, 1\}$: whether to detect spare parts 1 (0 is not detected, 1 is detected)
 $S_2 \in \{0, 1\}$: whether to detect spare parts 2 (0 is not detected, 1 is detected)
- 2) Stage 2: Finished product inspection decision; state representation: S_3
 $S_3 \in \{0, 1\}$: whether to detect the finished product (0 is not detect, 1 is detect)

3) Stage 3: Decision making on the dismantling of substandard finished products (substandard products after testing and substandard products that have not been tested and have been brought to the market): state representation: S_4

$S_4 \in \{0, 1\}$: whether to disassemble the finished product that fails the test (0 is not disassembled, 1 is disassembled).

4) Stage 4: Decision-making for part 1 and part 2 after dismantling of substandard products; Status representation: (S_5, S_6)

$S_5 \in \{0, 1\}$: whether to detect Part 1; $S_6 \in \{0, 1\}$: whether to detect Part 2.

The 16 actions generated were sequentially coded as shown in Table 3:

Table 3. 16 Action Codes

Index	Finished Product Inspection	Defective Product Disassembled	Component1 Inspected After Disassembly	Component2 Inspected After Disassembly
1	0	0	0	0
2	0	0	0	1
...
16	1	1	1	1

2.3. Modeling and solving

2.3.1. Economic benefit assessment model.

The core objective of the optimal decision-making scheme is to achieve cost minimization. In the enterprise production process, each stage involves a certain amount. Therefore, of cost inputs, including procurement costs, quality testing costs and so on, it is necessary to make an optimal decision on the cost inputs of each production stage through the analysis of economic efficiency assessment in order to achieve the goal of maximizing the economic efficiency of the enterprise [6]. The total cost components are as follows:

$$\text{Total Cost} = C_{\text{inspect}} + C_{\text{assemble}} + C_{\text{disassemble}} + C_{\text{replace}} + C_{\text{disposal}} \quad (1)$$

C_{inspect} : the cost required for quality testing of spare parts or finished products;

C_{assemble} : the cost required to assemble as spare part into a finished product;

$C_{\text{disassemble}}$: the costs incurred in dismantling sub-standard products;

C_{replace} : losses incurred on the exchange of finished products that have reached the market;

C_{disposal} : the cost of discarding parts or finished products that fail testing;

2.3.2. Optimization strategy: dynamic programming.

Dynamic programming [7] allows solving the optimal solution of the whole problem contains the optimal solutions of its sub problems, and by saving the solutions of the sub-problems (memorization), repeated calculations can be avoided. Since the optimal decision scheme of this problem is composed of the optimal solutions of each stage, and the decision of each stage will affect the decision of the next stage, the dynamic programming method is chosen to solve it.

Define 4 stage states and 16 behaviors: there are 2 spare parts, each of which can be chosen to be detected or not, for a total of $2^2=4$ ways of combining strategies. Thus, the decision scheme containing the four stages has a total of $4 \times 2^4=64$ strategy combinations

1) Decision-making phase of spare parts testing

The enterprise needs to decide whether or not to test spare part 1 and spare part 2. If they are tested, they incur certain testing costs and the detected unqualified products are discarded; if they are not tested, the spare parts go into the assembly of the finished product.

The cost of testing parts 1, 2 : $V(1,1,1) = C_{inspect1}$, $V(1,2,1) = C_{inspect2}$;

2) Finished product testing decision making stage

The assembled finished product contains qualified spare parts and unqualified spare parts, while the finished product composed of unqualified spare parts must be unqualified. and the finished product composed of qualified spare parts must not be qualified. Therefore, enterprises need to decide whether to test the assembled finished products. First, count the number of assembled finished products:

$$m = \begin{cases} 1 - p_1 & , \text{if part1 inspect, pat 2 not inspect} \\ 1 - p_2 & , \text{if part1 not inspect, pat 2 inspect} \\ 1 & , \text{if part1 pat 2 all not inspect} \\ \min\{1 - p_1, 1 - p_2\} & , \text{if part1 pat 2 all inspect} \end{cases} \quad (2)$$

p_1 denotes the defective rate of spare parts 1, p_2 denotes the defective rate of spare parts 2

Next, calculate the cost of finished assembly:

$$C_{Assemble} = m \times C_{assemble} \quad (3)$$

Next, calculate the cost of testing the finished product:

$$V(2,1) = m \times C_{assemble} + m \times C_{inspect} \quad (4)$$

3) Decision making stage for dismantling substandard finished products

For detected unqualified finished products and untested finished products flowing into the market, enterprises will be discarded for treatment, or sent back to the supplier for dismantling and processing, recovering spare parts for reuse and reducing production losses, but dismantling will incur certain costs. Therefore, enterprises need to decide whether or not to dismantle unqualified finished products

If disassembled, the cost incurred is:

$$V(3,1) = m \times p \times C_{disassemble} \quad (5)$$

After dismantling the substandard finished products need to re-enter the spare parts testing decision stage and finished product testing decision stage

If not disassembled, the finished product will be discarded and the cost of purchasing spare parts land 2 becomes a sunk cost:

$$V(3,0) = m \times p \times (C_1 + C_2) \quad (6)$$

p denotes the defective rate of the finished product, C_1 and C_2 denote the purchase unit price of spare part 1 and spare part 2 respectively

4) Decision making stage for inspection of dismantled parts

After the disassembly of non-conforming finished products part 1 and part 2 re-enter the parts testing decision measurement stage, you can choose whether to carry out testing or not, the recovery of non-conforming products will re-enter the non-conforming finished products disassembly decision stage. The formula for the testing cost is as follows:

The cost of testing parts 1,2 is

$$V(4,1,1) = C_{inspect1}, V(4,2,1) = C_{inspect2} \quad (7)$$

2.3.3. Model results and analysis.

The costs corresponding to the 64 combination strategies are calculated by MATLAB, and the box line diagram of the cost distribution of the plotted strategies is shown in Figure1, from which it can be seen that under the same data of defective rate, inspection cost, and swapping loss, the costs generated by different inspection combination strategies vary greatly. The maximum cost in the inspection scheme reaches about \$27.4, and the minimum is only about \$7.4(the calculation process assumes that the number of parts is the number of units 1). Therefore, choosing a good inspection scheme plays an important role in reducing costs and maximizing profits.

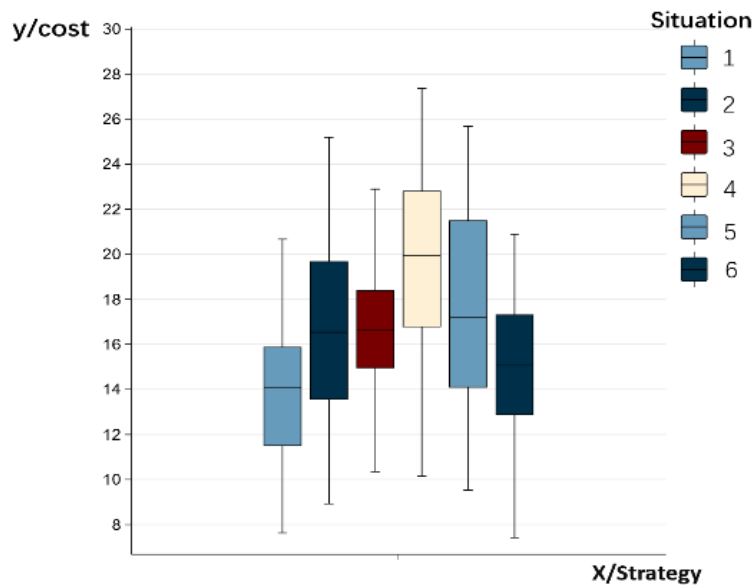


Figure 1. Boxplot of Cost Distribution

After the analysis of the four stages of the production process for dynamic planning to optimize the decision path and economic efficiency assessment, the optimal detection decision scheme and the corresponding cost including the decision and scheme of each stage in six different cases are obtained, as shown in Table 4:

Table 4. Optimal Decision Solution

Situation	Part 1 Inspect	Prat 2 Inspect	Final Product Inspected	Defective Product Disassembly	Part 1 Inspect After	Part 2 Inspect After	Cost (\$)
1	No	No	No	No	No	No	7.62
2	No	No	No	No	No	No	8.92
3	No	No	Yes	Yes	No	No	10.35
4	Yes	No	Yes	Yes	No	No	10.44
5	No	No	No	No	No	No	9.52
6	No	No	No	No	No	No	7.426

From the above table, it can be seen that in the parts part (part 1 and part 2), enterprises choose not to carry out testing, which indicates that the defective rate is low or the cost of testing is high, so skipping the testing can save costs. In the finished product part, such as scenarios 3and 4, when the exchange loss is at a high level and the defective rate is not low, the enterprise needs to strictly control the quality of the finished product and test the finished product to minimize the impact of the exchange

loss on the overall profit. Take Program 4 as an example, the defective rate of spare parts 1 is relatively high and the detection cost is low, the defective rate of finished products is high and the detection cost is relatively low, while the exchange loss is at a high level, the enterprise should choose not only to test the spare parts 1, but also to test the finished products, to avoid the unqualified finished products into the market to produce a greater cost of investment.

2.4. Sensitivity analysis

There are a large number of environmental variables that directly affect profit in the optimal decision scheme problem, including other variables such as the unit price of parts purchased, the cost of assembly of finished products, and the cost of testing parts or finished products. A sensitivity analysis is performed [8], so that these variables fluctuate continuously up and down by 5% from their initial values, and a schematic diagram is plotted between each of these three parameters and the total cost:

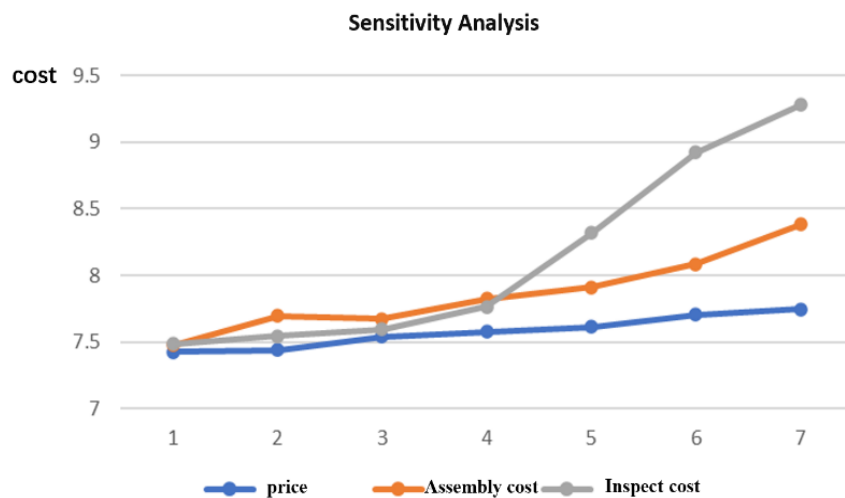


Figure 2. Sensitivity Analysis

As can be seen from the above figure 2, the dynamic planning model is more sensitive to the 2 parameters of testing cost and purchase unit price, in which the change of testing cost of the parts will have a greater impact on the total cost of the program, and when the cost of testing is getting bigger and bigger, the total cost of the program 2 grows at a faster rate. Therefore, extra attention should be paid to their values in the actual participation in business decision-making.

3. Markov decision process

3.1. Problem assumptions

The problem is further complicated by the fact that Figure 3 gives 2 processes and 8 spare parts, and Table 5 gives the specific data:

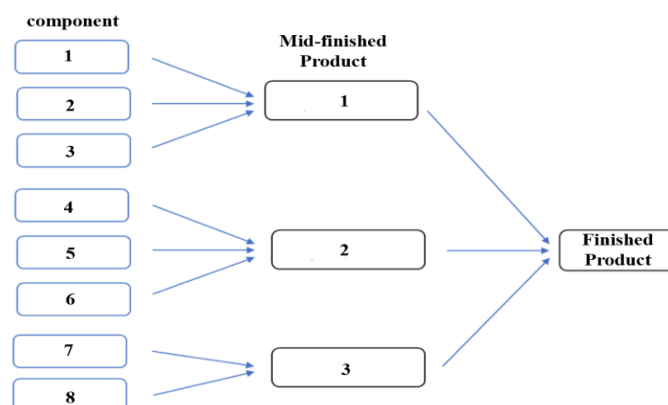


Figure 3. Process

Table 5. Production Conditions

component	defect rate	Inspection cost (\$)	Mid-finished Product	defect rate	Inspection cost	Assembly cost (\$)	Disassembly Cost (\$)
1	10%	1	1	10%	4	8	6
2	10%	1	2	10%	4	8	6
3	10%	2	3	10%	4	8	6
4	10%	1					
5	10%	1	Finished	10%	6	10	8
6	10%	2					
7	10%	1	Market price			Replacement loss	
8	10%	2	Finished	200		40	

3.2. Modeling

1) State definition and discretization

- Parts detection stage: $s_{1i} \in \{0,1\}$ Whether to detect parts, $i = 1, 2, \dots, 8$
- Semi-finished product detection disassembly stage: $s_{2i} \in \{0,1\}$ Whether to detect semi-finished products $d_{2i} \in \{0,1\}$ Whether to disassemble the semi-finished product
- Finished product detection and disassembly stage: $s_3 \in \{0,1\}$ Whether to detect the finished product or not $d_3 \in \{0,1\}$ Whether to disassemble the finished product.
- The stage of dismantling of substandard finished products flowing into the market: $d_4 \in \{0,1\}$ No dismantling of non-conforming finished products that have not been detected as having reached the market

2) Action definition and discretization

Define the corresponding decision action for each state and enumerate all possible actions, the action space may change with the state

3) Defining state transfer probabilities

Determine, for each state, the probability of state transfer after performing a specific action, based on the rate of inferiority $P(s' | s, a)$

4) Defining the reward function

Given the current state and decision, define the immediate reward function $R(s, a)$ The reward function should reflect the cost objective of the organization:

$$R(s, a) = -(C_{inspect} + C_{assemble} + C_{disassemble} + C_{disposal} + C_{replace}) \quad (8)$$

5) Solving the optimal policy

The central goal of solving the MDP [9,10] is to minimize the long-term cost by choosing the optimal action “a” at each state “s”. The Bellman equation for the Markov decision process [4] is:

$$V(s) = \max_a [R(s, a) + \gamma \sum_s P(s | s, a) V(s')] \quad (9)$$

$V(s)$ is a function of the value of state “s”, denoting the maximum expected gain from state “s”; γ is a discount factor that reflects the impact of future earnings on current decisions.

3.3. Model solving

Calculate the total number of combinations of strategies, stage 1: parts 1-8 can choose whether to detect, a total of $2^8=256$ kinds of combinations; stage 2: set the non-detection of the finished product must not be disassembled, while deleting the duplicate results, a total of 27 kinds of combinations of strategies stage 3: after the detection of unqualified finished product dismantling, after the detection of the finished product does not disassemble with the non-detected unqualified finished product directly into the market, a total of 3 kinds of cases; stage 4: without detection Unqualified finished products flowing into the market can choose whether to disassemble or not, a total of $2^1=2$ combinations. The total strategies are $256 \times 27 \times 3 \times 2 = 41472$ different strategy combinations.

1) Decision-making phase of spare parts testing

The firm decides whether or not to test the spare parts, and if it does, it incurs a testing cost, and those that fail the test will be discarded, so the loss of discard is a purchase cost; if it does not test, it incurs no cost. The reward function is expressed as:

$$R(S_1, a) = \sum_{i=1}^8 (S_{1i} \times C_{inspect}) + \sum_{i=1}^8 (S_{1i} \times P_{1i} \times C_i) \quad (10)$$

P_{1i} Indicates the part of the defective parts

The process quantity for calculating the probability of state transfer in stage 1 is:

$$P_1(S' | S, a) = \begin{cases} 1 - P_{1i}, & \text{if } S_{1i} = 1 \\ 1, & \text{if } S_{1i} = 0 \end{cases} \quad (11)$$

Since it is assumed that the number of each part required to make up the semi-finished product is equal, the number of parts with the least number of parts is chosen as the number of semi-finished products to be assembled, and the state transfer probability is expressed as:

$$P_1'(S' | S, a) = \begin{cases} 1, & \text{if part } i \text{ not inspect} \\ \min \{1 - P_{1i}, 1\}, & \text{if parts must be tested at least once} \end{cases} \quad (12)$$

2) Semi-finished product testing and dismantling decision-making stage

In this stage, assembly costs are incurred when parts are assembled into a certain number of semi-finished products. At the same time, the enterprise needs to decide whether to test or disassemble the semi-finished products. If testing, a certain amount of testing costs will be incurred, for the detection of substandard semifinished products need to make further decisions on whether to dismantle the treatment. If disassembled disassembly cost will be incurred: if not disassembled for disposal, disposal loss will be incurred. The reward function is expressed as:

$$R(S_2, a, d) = \sum_{j=1}^3 (S_{2j} \times C_{inspect}) + \sum_{j=1}^3 (S_{2j} \times d_{2j} \times \sum_{i \in j} C_i \times T_j) \\ + \sum_{j=1}^3 (S_{2j} \times (1 - d_{2j}) \times C_{disassemble} \times T_j) + P_1' \times C_{assemble} \quad (13)$$

T_j Indicates the actual rate of defective semi-finished products:

$$T_j = \prod_{i \in j} (1 - S_{1i}) \times P_1 + P_2 (1 - P_1)^{\sum_{i \in j} (1 - S_{1i})} \quad (14)$$

P_1 denotes the defective rate of spare parts, and P_2 denotes the defective rate of semi-finished products.

The process quantity for calculating the probability of state transfer in stage 2 is:

$$P_2(S' | S, a) = \begin{cases} 1 - T_j, & \text{if } S_{2j} = 1 \\ 1, & \text{if } S_{2j} = 0 \end{cases} \quad (15)$$

By calculating the process volume, the state transfer probability for stage 2 can be expressed as:

$$P_2'(S' | S, a) = \min\{P_2\} \quad (16)$$

3) Finished product testing and dismantling decision-making phase

The cost of this phase consists of finished product assembly cost, finished product inspection cost, loss of disposal of non-conforming finished product after inspection and dismantling cost of nonconforming finished product after inspection. The reward function is expressed as:

$$R(S_3, a, d) = TC_{inspect} + TC_{disposal} + TC_{disassemble} + P_2' \times C_{assemble} \quad (17)$$

TC represents the cost associated with the finished product, and the actual finished product defective rate is expressed as:

$$T_2 = \prod_{j=1} (1 - S_{2j}) \times T_1 + P_3(1 - T_1)^{\sum_{j=1} (1-S_{2j})} \quad (18)$$

The state transfer probability for stage 3 is denoted as:

$$P_3'(S' | S, a) = \begin{cases} 0, & \text{if } S_{3i} = 1 \\ T_2, & \text{if } S_{3i} = 0 \end{cases} \quad (19)$$

4) Decision making phase of dismantling substandard finished products that have entered the market without testing

$$R(S_4, d) = T_2 \times C_{inspect} + T_2 \times C_{disassemble} + T_2 \times \sum_{i=1}^8 C_i \quad (20)$$

3.4. Analysis of model results

The data for each indicator was substituted into MATLAB and MDP was used to calculate two similar decision options:

Option 1: For the parts part, test part 1 and do not test parts 2-8, for the semi-finished part, do not test and do not disassemble semi-finished products 1-3, for the finished part, do not test and do not disassemble finished products, and disassemble and dispose of unqualified finished products that flow into the market Under this decision scheme, the total cost of production for the firm is \$32.8023 (assuming the number of parts is unit 1).

Option 2: For the parts section, parts 1-8 are not tested, for the semi-finished parts section, semi-finished products 1-3 need to be tested, for the finished parts section, finished products need to be tested, and unqualified finished products need to be disassembled. Under this decision scheme, the total cost of production for the firm is \$32.8065 (assuming the number of parts is unit 1).

Comparison reveals that Option 2 tested semi-finished products and costs, while Option 1 did not test semi-finished products and finished products, and it can be inferred that the exchange loss of Option 2 is smaller than the exchange loss of Option 1. Considering that in the case of similar costs, the enterprise reputation caused by the flow of unqualified finished products into the market, fluctuating logistics costs and other switching losses for the enterprise's future production and profitability of a potential threat. Therefore Option 2 is chosen as the optimal decision-making program.

As can be seen from scheme 2, in the parts part, the enterprise should choose not to test the parts, which indicates that the defective rate of the parts is low and skipping the test can save the cost. And in the semi-finished part, the enterprise chooses to test the semi-finished product, which indicates that in the semi-finished product stage, the enterprise needs to strictly control the quality in order to avoid the larger exchange loss caused by the unqualified finished product flowing into the market. For the finished product part, the enterprise should test the finished product and disassemble the unqualified products, and the disassembled parts re-enter the parts testing decision making stage, which indicates that the enterprise needs to strictly control the final stage of production and fully recycle the disassembled spare parts in order to reduce the production loss and save costs.

4. Conclusion

This paper addresses the decision-making problem in factory production, firstly, by combining the model of dynamic planning and economic benefit assessment, calculating the optimal strategy in basic production, and then further complex production process. Combined with Markov decision process model, optimize the benefit assessment mechanism, and finally through the sensitivity analysis, fully demonstrates that the model in this paper can efficiently, accurately and flexibly provide decision-making advice for the enterprise production. which improves the real problems of resource waste and low efficiency to a certain extent, and promotes the development of productivity.

With the development of electronic products, there may be more stages in the production process in the future that require decision-making, and more complex influencing variables (such as environmental factors, equipment failure rates, raw material price fluctuations, etc.) will be introduced. Therefore, future research can take more factors into account based on the model in this paper, combined with the multi-objective optimization model, so that production decisions can not only focus on cost, but also comprehensively consider key indicators such as production efficiency and product quality, to further improve the accuracy, flexibility and sustainability of production decisions.

References

- [1] Zhang R, Wei F. A Dynamic Stochastic Production Decision Model Based on Analytical-Simulation Feedback [J]. *Computer Integrated Manufacturing Systems*, 2005, 11 (12): 0.
- [2] Yang C. Design and Analysis of Dynamic Programming Algorithm for Multi-stage Flow Shop Scheduling [J]. *Computer Programming Techniques and Maintenance*, 2023, 11 (11): 20 - 22+39.
- [3] Cai M. Solving the Optimal Production Strategy Problem Using Dynamic Programming and Lingo [J]. *Forest Area Teaching*, 2015, (11): 80 - 81.
- [4] Marfuah, Umi. Dynamic programming approach in aggregate production planning model under uncertainty [J]. *International Journal of Advanced Computer Science and Applications*, 2023, 14 (3).
- [5] Yan Y, Shuai P, Zhao Y. 0-1 Programming Mathematical Model for Automotive Spare Parts Production Scheduling [J]. *Industry and Technology Forum*, 2022, 21 (11): 40 - 42.
- [6] Li Q., Energy-saving and Environmental Protection Benefits Evaluation of Intelligent Unmanned Factories in New Energy Production [J]. *Environment and Development*, 2024, 6 (2).
- [7] Forootani, Ali, et al. A stochastic dynamic programming approach for the machine replacement problem[J]. *Engineering Applications of Artificial Intelligence*, 2023, 118: 105638.
- [8] Wang X, et al. Sensitivity Analysis and Adaptability Study of Rice SM Rice Model Parameters [J]. *Journal Of Agricultural Big Data*, 2023, 5 (2): 97 - 108.
- [9] Wu J, Zhao J, Sun C., et al. Predicting Resource Allocation: Unsupervised Learning of Markov Decision Process [J]. *Science China: Information Sciences*, 2024, 54 (08): 1983 - 2000.
- [10] Wu Q. Order scheduling optimization in manufacturing enterprises based on MDP and dynamic programming [J]. *Scientific Reports*, 2023, 13 (1): 9783.