

➤ Personal Introduction



RED BIRD CHALLENGE CAMP

PORTFOLIO FOR HKUST(GZ)

YIXIAO ZHANG

📍 Chengdu, China

☎ +86 18202429280

✉ q18202429280@outlook.com

🌐 [linkedin.com/in/yixiao-zhang-2042ab2a3/](https://www.linkedin.com/in/yixiao-zhang-2042ab2a3/)

💻 github.com/13558882230

Northeastern University
Bachelor of Management
Major: Information Management & Information Systems
Sep 2022 - Jul 2026 | GPA: 3.47/5 (Top 25%)



Hanyang University
Hanyang International Winter School
GPA: 4.25/4.5 (Top 5%)



Data Analysis & Visualization: Proficient in Python, R, Jupyter Notebook for data mining; skilled in Tableau and Streamlit for interactive visualizations.

Full-Stack Development: Backend: Flask framework for stable service development; Frontend: Angular, HTML/CSS/JavaScript for responsive interface design.

Database Management: MySQL (relational) and Elasticsearch (vectorized) .

Tools & Infrastructure: Git (version control), VS Code, Linux, Docker (containerization).

Programming Languages: Python, Java, R, SQL.

Language Proficiency: CET-6 (English); Mandarin Level 2-B.

HONORS AND AWARDS

Business Elite Challenge Accounting and Business Case Competition(**National First Prize**)

Northeastern University 11th Overseas Economic Management Scholars Seminar(**Outstanding Camper**)

Estonian National Summer School Full Scholarship (2025) (17/500)

€1,050 scholarship by Estonian Education and Youth Board.

Bucharest Summer University Partial Scholarship (2025)

Selected for partial funding at 19th BSU (750+ participants).



➤ Work experience

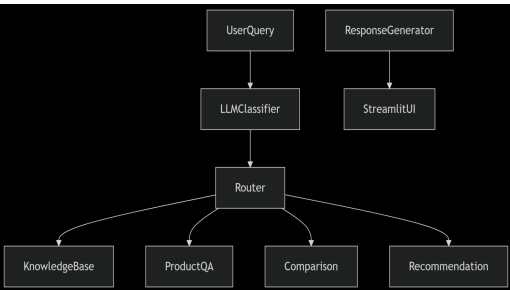
Xiaoduo Tech (7-9 2024) Prompt Engineer



1. System Architecture

E-commerce query routing system using RAG architecture

- Developed LLM-based intent classifier
- Built lightweight Elasticsearch vector database
- Created Streamlit demo for AI interaction showcase



2. Performance Metrics

Metric	Before	After
Latency	450ms	280ms
Accuracy	72%	85%
Error Rate	5.1%	0.7%

3. Business Impact

10,000+ daily queries →
85% accurate routing →
34% higher conversions
→ \$28K monthly revenue increase

Algorithm 1 Streamlit Product QA System

```
Initialization
1: Connect to Elasticsearch: es ← Elasticsearch(host = localhost)
2: Initialize session: session.state ← {}
Main Application Flow
3: procedure MAIN
4:   Display sidebar navigation with 5 pages
5:   selected_page ← st.sidebar.radio()
6:   if selected_page == "" then
7:     QUESTION REWRITING
8:   else if selected_page == "" then
9:     PRODUCT QA
10:  else if selected_page == "" then
11:    KNOWLEDGE BASE
12:  else if selected_page == "" then
13:    PRODUCT COMPARISON
14:  else
15:    PRODUCT RECOMMENDATION
16:  end if
17: end procedure
Core Functions
18: function QUESTION REWRITING
19:   user_input ← st.text_input()
20:   if user_input ≠ "" then
21:     results ← es.search(index = product_names, query = user_input)
22:     prompt ← PROMPT_FMT1(results, user_input)
23:     response ← OpenAI(prompt)
24:     session.state.questions ← parse(response)
25:   end if
26: end function
27: function PRODUCT QA
28:   if session.state.questions ≠ "" then
29:     selected ← st.select_box(options = session.state.questions)
30:     params ← es.get_params(selected.product_name)
31:     prompt ← PROMPT_FMT2(selected.question, params)
32:     st.write(OpenAI(prompt))
33:   end if
34: end function
35: function PRODUCT COMPARISON
36:   if selected then
37:     selected ← choose.question.from.collection()
38:     params_list ← [es.get_params(name) ∀ name ∈ selected.product_names]
39:     prompt ← PROMPT_FMT3(selected.question, params_list)
40:     st.write(OpenAI(prompt))
41:   end if
42: end function
Key Data Structures
43: PROMPT_FMT1 ← XML template for question rewriting
44: PROMPT_FMT2 ← Template for single product QA
45: PROMPT_FMT3 ← Template for product comparison
46: session.state ← {questions : List[Dict]}
```



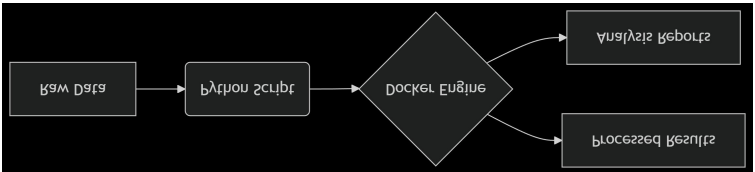
Bairen BioTech (10-12 2024) Platform Development



1. Project Overview

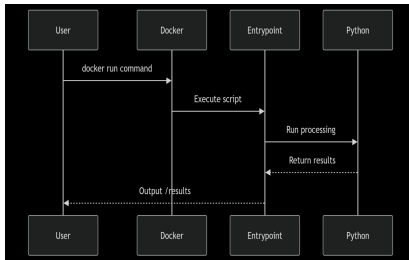
One-command bioinformatics analysis pipeline

- Developed Python data processing scripts
- Containerized environment with Docker
- Created entrypoint automation system



2. Workflow:

User runs docker run -v data:/input bioren-image
docker_entrypoint.sh executes:



- Data validation
- Pipeline sequencing
- Result export

3. Performance Gains

Metric	Manual	Automated
Setup Time	3 hrs	5 min
Analysis Speed	8 hrs	2.5 hrs
Reproducibility	Low	100%

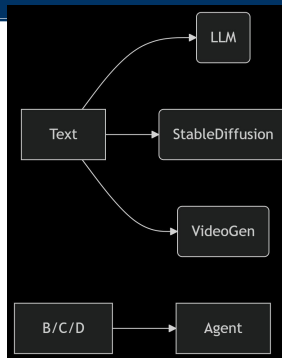
Migu Music(12.2024-3.2025) AI Data Intern



1. Project Overview

Spring Festival AI Agent Development

- Optimized prompts for LLM/text-to-image/text-to-video
- Built structured label recognition system
- Created 200+ festival-themed assets

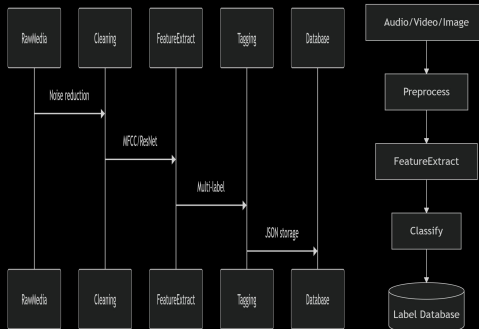


2. Prompt Engineering Optimization Approach:

```
# Before
prompt = "Spring Festival image"

# After (Optimized)
prompt = """
Chinese New Year scene:
- Red lanterns hanging
- Gold ingots decoration
- Tiger motif (Year of Tiger)
- Style: Traditional Chinese painting
- Aspect ratio: 16:9
"""
```

3. Structured Label Recognition



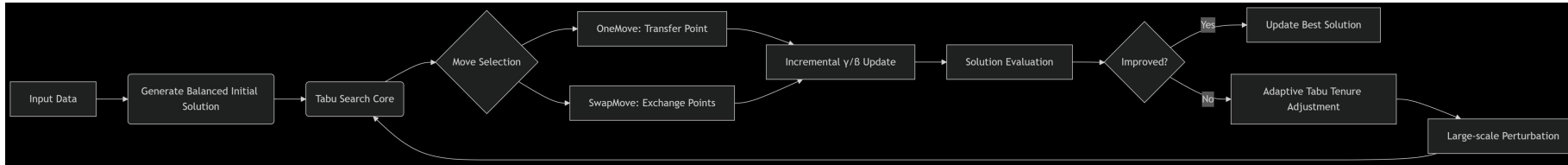
4. Performance Metrics

Metric	Before	After
Asset Production	2/hr	8/hr
Label Accuracy	78%	92%
Rejection Rate	40%	12%

➤ Research results

Balanced Minimum Sum-of-Squares Clustering Problem

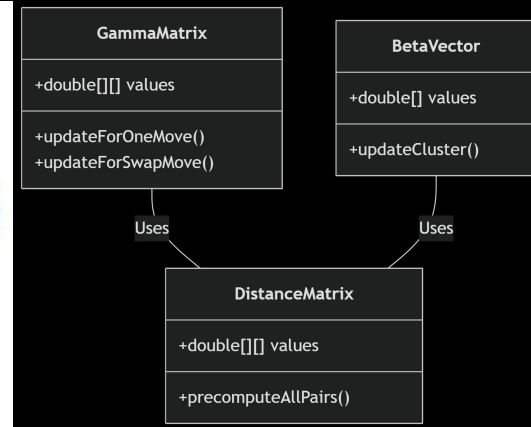
National Natural Science Foundation Project | Java Implementation



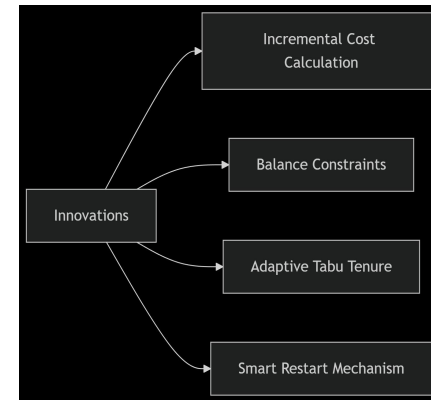
1.Core Technical Innovations

Incremental Evaluation System:

$$\Delta_f(\text{OneMove}(p, C_g, C_h)) = \frac{\beta[g] - \gamma[p][g]}{|C_g| - 1} + \frac{\gamma[p][h] + \beta[h]}{|C_h| + 1} - \left(\frac{\beta[g]}{|C_g|} + \frac{\beta[h]}{|C_h|} \right)$$
$$\Delta_f(\text{SwapMove}(p, q)) = \frac{\gamma[q][g] - \gamma[p][g] - \|p - q\|^2}{|C_g|} + \frac{\gamma[p][h] - \gamma[q][h] - \|p - q\|^2}{|C_h|}$$



2.Innovation Highlights



Algorithm 1 Tabu Search for Balanced Clustering

```
1: Input:
2:   Data points  $D = \{x_1, \dots, x_n\}$ 
3:   Number of clusters  $k$ 
4:   Maximum runtime  $T_{max}$ 
5: Output: Balanced clusters  $C = \{C_1, \dots, C_k\}$  with minimal SSE
6: procedure TABUSEARCH( $D, k, T_{max}$ )
7:   Initialize distance matrix  $dist[i][j] \leftarrow \|x_i - x_j\|^2, \forall i, j$ 
8:   Generate initial balanced solution  $C \leftarrow$ 
   generateBalancedKMeansSolution()
9:    $bestCost \leftarrow$  calculateCost( $C$ )
10:   $globalBest \leftarrow C, globalBestCost \leftarrow bestCost$ 
11:  Initialize tabu list  $lastMoveIteration[n][k] \leftarrow 0$ 
12:  Initialize  $\gamma[p][q] \leftarrow \sum_{q \in C_g} dist[p][q]$ 
13:  Initialize  $\beta[g] \leftarrow \sum_{p \in C_g} dist[p][q]$ 
14:   $tabuTenure \leftarrow \sqrt{n}/k$ 
15:  while time  $< T_{max}$  and restart count  $< maxRestarts$  do
16:     $noImprovement \leftarrow 0$ 
17:    while time  $< T_{max}$  and  $noImprovement < maxNoImprovement$ 
18:      do
19:        Find best non-tabu move (OneMove or SwapMove):
20:        OneMove:  $\Delta f = \frac{\beta[g] - \gamma[p][g]}{|C_g| - 1} + \frac{\gamma[p][h] + \beta[h]}{|C_h| + 1} - \left( \frac{\beta[g]}{|C_g|} + \frac{\beta[h]}{|C_h|} \right)$ 
21:        SwapMove:  $\Delta f = \frac{\gamma[q][g] - \gamma[p][g] - dist[p][q]}{|C_g|} + \frac{\gamma[p][h] - \gamma[q][h] - dist[p][q]}{|C_h|}$ 
22:        if move improves solution or accepted by SA criteria then
23:          Execute move
24:          Update  $\gamma$  and  $\beta$  incrementally
25:          Update tabu list with current iteration
26:           $bestCost \leftarrow bestCost + \Delta f$ 
27:          if  $bestCost < globalBestCost$  then
28:             $globalBest \leftarrow C, globalBestCost \leftarrow bestCost$ 
29:             $noImprovement \leftarrow 0$ 
30:          else
31:             $noImprovement \leftarrow noImprovement + 1$ 
32:          end if
33:          Adjust  $tabuTenure$  dynamically
34:        end while
35:        if no improvement or temperature too low then
36:          Restore  $C \leftarrow globalBest$ 
37:          Perform large-scale perturbation (swap 40% of points)
38:          Reinitialize  $\gamma$  and  $\beta$ 
39:          Increment restart count
40:        end if
41:      end while
42:    return  $globalBest$ 
43: end procedure
44: procedure GENERATEBALANCEDKMEANSOLUTION
45:   Select initial centers using max-min distance
46:   Calculate target cluster sizes:  $base \leftarrow \lfloor n/k \rfloor, remainder \leftarrow n \bmod k$ 
47:   Assign points to nearest center while respecting size constraints
48:   return balanced clusters
49: end procedure
```

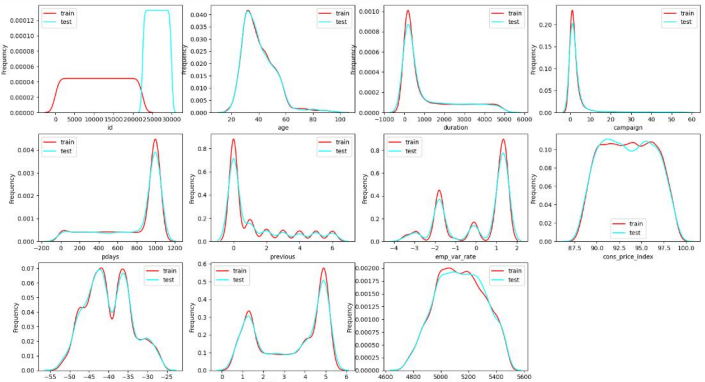
3.Comparative Analysis

Algorithm	Optimal Solutions	Avg Dev (%)	Speed
bk-means	30/160	2.09	1.0x
VNS-LIMA	47/160	1.38	1.2x
Ours	159/160	0.49	1.4x

- Developed an enhanced K-means initialization protocol (20+ iterations) to generate high-quality initial solutions
- Engineered a real-time exchange cost evaluation system using γ/β matrices enabling $O(1)$ operation cost calculation
- Implemented Tabu Search with optimized OneMove/SwapMove operations, improving solution quality by 25%
- Integrated population-based algorithms for multi-path exploration, increasing global search efficiency by 40%

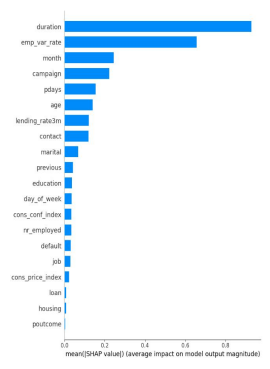
Main Academic Achievements

1.Data Distribution



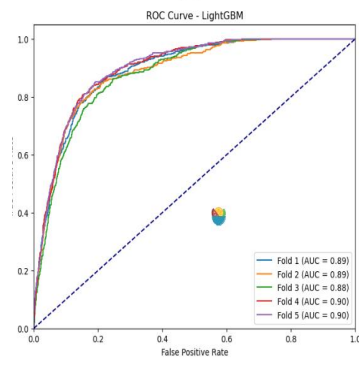
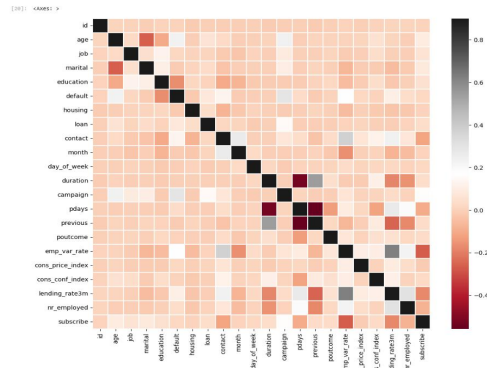
2.LightGBM Modeling Summary:

- ▶ Algorithm: Gradient boosting with tree-based learning, regularization against overfitting
- ▶ Performance: 5-fold CV shows strong AUC (0.88-0.90), confirming high predictive accuracy
- ▶ Key Drivers: SHAP analysis identifies duration and emp_var_rate as top features, with *age/lending_rate3m* less impactful

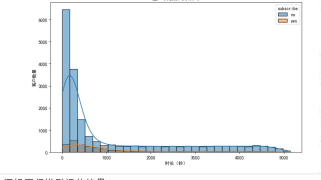
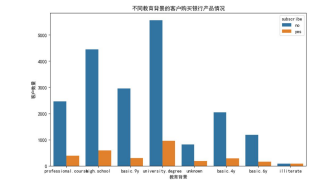
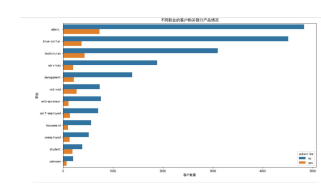
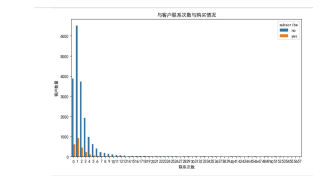
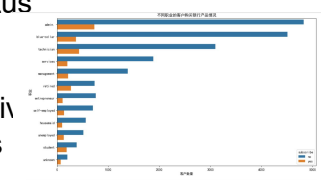
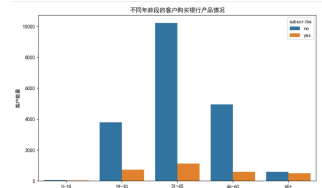


3.Key Insights from Correlation Heatmap:

- ▶ Unexpected Pattern: High id-age correlation (potential data artifact)
- ▶ Socioeconomic Drivers: Moderate positive links (job/marital/education → subscription)
- ▶ Financial Factors: Housing/loan status boosts subscriptions
- ▶ Economic Signals: Negative employment_var_rate impact vs. positive *lending_rate3m/nr_employed* effects



4.Customer Profile Analysis



Model Comparison Highlights:

- ▶ Random Forest: 88.27% accuracy (top performer), excels in minority-class detection
 - ▶ XGBoost: 88.06% (slightly lower despite handling complexity well)
 - ▶ Logistic Regression: 86.46% (stable baseline)
- Conclusion: Random Forest preferred for balanced precision/recall/F1 and interpretability.

逻辑回归模型评估结果:

	precision	recall	f1-score	support
no	0.88	0.97	0.93	5860
yes	0.46	0.16	0.23	890
accuracy				
macro avg	0.67	0.56	0.58	6750
weighted avg	0.83	0.86	0.83	6750

Accuracy: 0.8645925925925926

随机森林模型评估结果:

	precision	recall	f1-score	support
no	0.90	0.98	0.94	5860
yes	0.63	0.27	0.37	890
accuracy				
macro avg	0.76	0.62	0.65	6750
weighted avg	0.86	0.88	0.86	6750

Accuracy: 0.8826666666666667

XGBoost模型评估结果:

	precision	recall	f1-score	support
no	0.91	0.95	0.93	5860
yes	0.57	0.40	0.47	890
accuracy				
macro avg	0.74	0.68	0.70	6750
weighted avg	0.87	0.88	0.87	6750

Accuracy: 0.8809259259259259

Algorithm 1 Bank Customer Purchase Behavior Analysis

- 1: **Input:** Training dataset (train.csv), Testing dataset (test.csv)
- 2: **Output:** Analysis results of customer features vs purchase behavior
- 3: **Step 1: Data Loading**
- 4: Load training data: `train_data ← pd.read_csv("d:/Desktop/20221188/train.csv")`
- 5: Load testing data: `test_data ← pd.read_csv("d:/Desktop/20221188/test.csv")`
- 6: **Step 2: Data Exploration**
- 7: Display training data preview: `print(train_data.head())`
- 8: Display testing data preview: `print(test_data.head())`
- 9: Check missing values: `print(train_data.isnull().sum())`
- 10: Describe age distribution: `print(train_data['age'].describe())`
- 11: **Step 3: Data Preprocessing**
- 12: Fill missing values (forward fill): `train_data.fill(inplace = True)`
- 13: `test_data.fill(inplace = True)`
- 14: Check age outliers:
- 15: Define bounds: `age_lower ← 0; age_upper ← 120`
- 16: Filter outliers: `print(train_data[(train_data['age'] > age_upper) | (train_data['age'] < age_lower)])`
- 17: **Step 4: Feature Engineering (Age Grouping)**
- 18: Define bins: `bins ← [0, 18, 30, 45, 60, 100]`
- 19: Define labels: `labels ← ['0 - 18', '19 - 30', '31 - 45', '46 - 60', '60 +']`
- 20: Create age groups:
- 21: `train_data['age-group'] ← pd.cut(train_data['age'], bins = bins, labels = labels)`
- 22: `test_data['age-group'] ← pd.cut(test_data['age'], bins = bins, labels = labels)`
- 23: **Step 5: Visualization**
- 24: Plot age group vs purchase: `plt.figure(figsize = (10, 6))`
- 25: `sns.countplot(x = 'age-group', hue = 'subscribe', data = train_data)`
- 26: `plt.title("Customer Purchase by Age Group")`
- 27: `plt.show()`
- 28: Plot job type vs purchase: `plt.figure(figsize = (14, 8))`
- 29: `order ← train_data['job'].value_counts().index`
- 30: `sns.countplot(y = 'job', hue = 'subscribe', data = train_data, order = order)`
- 31: `plt.title("Customer Purchase by Job Type")`
- 32: `plt.show()`
- 33: **Output:** Saved visualizations and analysis report

➤ Main Academic Achievements

Business Elite Challenge Accounting and Business Case Competition



Rebecca Group Strategic Business Analysis

1. Financial Snapshot

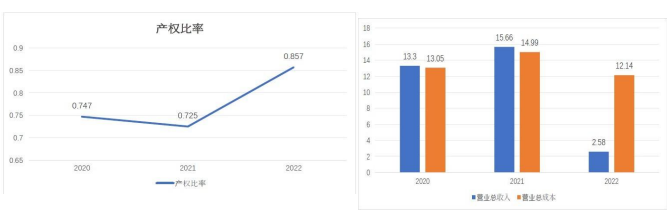
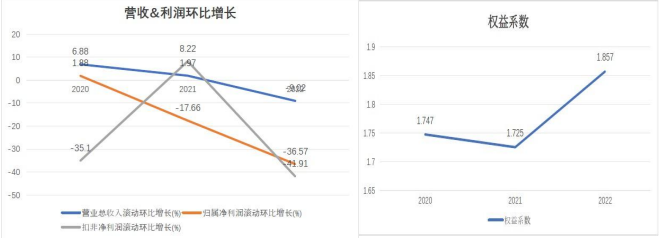
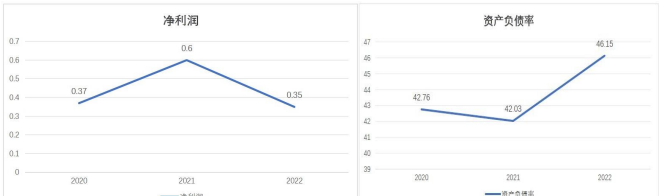
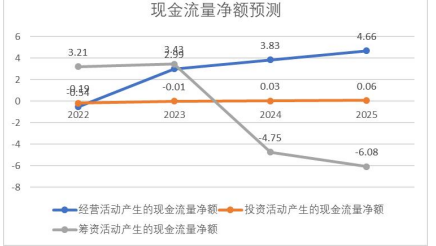
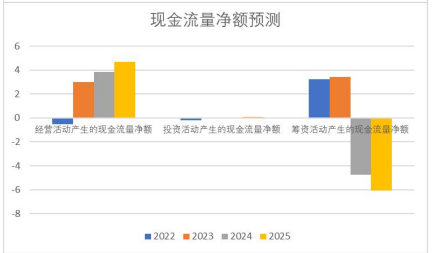
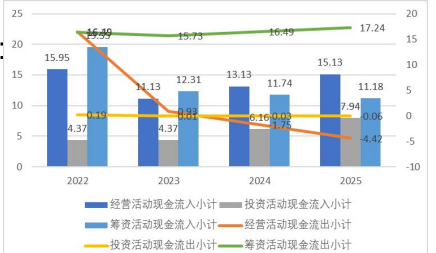
- Revenue: ¥1.33B → ¥1.26B (▼5.3%)
 - Overseas sales: 99% of total (Africa/Europe-focused)
- Profitability:
 - Net Profit: ¥38M → ¥34M (▼10.9%)
 - ROE: ▼ YoY (Operational inefficiency)
- Asset-Liability Shift:
 - ✓ Liquidity ▲: Cash reserves ↑13.88% (2022)
 - ✗ Debt risk: Short-term loans 57.9% (2022)

2. Core Strategic Frameworks

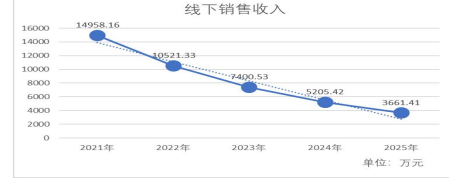
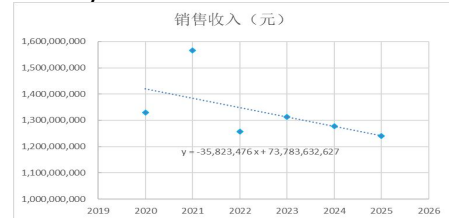
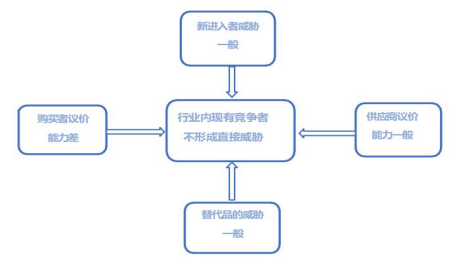
- PEST Analysis:
 - Political: RCEP export certification risks
 - Economic: Forex volatility (▼\$ impact)
 - Social: Customization demand ▲
 - Tech: E-commerce + AR try-on R&D
- Porter's 5 Forces:
 - Supplier Power: Low (Global sourcing)
 - Buyer Power: High (Omnichannel sales)
 - New Entrants: Moderate (E-commerce lower barrier)

3. Financial Forecast

Revenue: ¥13.1B→12.4B (▼5%) driven by online crash (▼65% to ¥0.33B) vs. overseas ¥12.0B (95%). 2025 costs spike: COGS ¥1.88B (▲147%), **expenses surge** (R&D▲133%, Sales▲242%, Mgmt▲167%). **Cash flow shifts:** Ops -¥54M→+¥466M; Financing +¥321M→-¥608M (debt wall). Risks: E-com failure, Africa costs, loan maturity.



瑞利卡商业画布				
关键合作伙伴	关键活动	价值主张	客户关系	客户细分
• 供应商	• 产品研发	• 高品质、时尚	• 专业售后	• 女性、发
• 分销商	• 品牌营销	• 舒适、耐用、独特	• 会员制度	• 发廊、沙龙、个人
• 品牌代言人/大使	• 渠道拓展	• 个性化搭配	• 新品优惠推送	
成本结构		核心资源	渠道进入	收入来源
• 原材料	• 设计师团队		• 实体店 (一二线城市)	• 发饰销售
• 运营	• 优质供应链		• 官网、电商平台	• 定制服务
• 营销	• 品牌影响力		• 社交媒体	• 附件销售

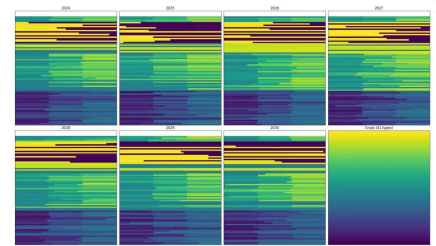
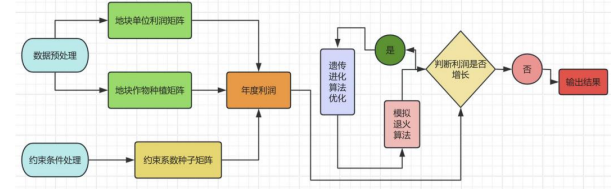


➤ Main Academic Achievements

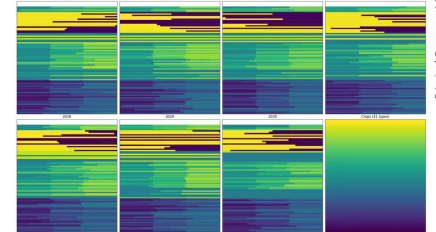
ChinaUndergraduate Mathematical Contest in Modelling Optimizing Crop Planting Strategies in Resource-Limited Mountainous Regions A Mathematical Modeling Approach for Sustainable Agriculture

PROBLEM 1: Deterministic Optimization
Challenge: Maximize profit (2024-2030) with stable market conditions

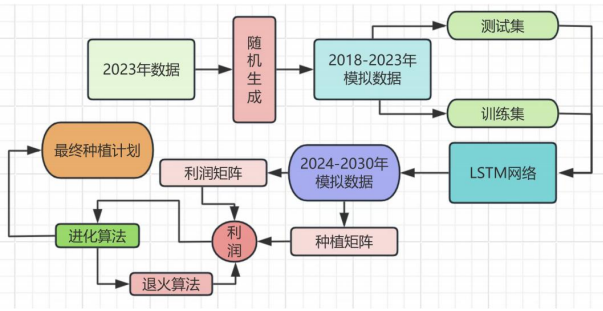
- Constraints:
- ▶ Land suitability (A-F classification)
 - ▶ Overproduction penalties (waste vs. 50% discount)
 - ▶ Crop rotation & continuous planting rules



(a) 情况 1 种植方案结果可视化

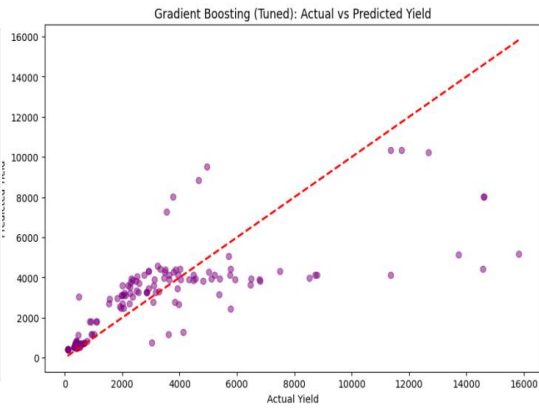
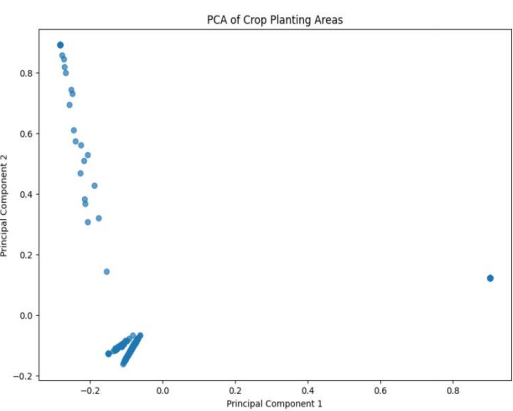


(b) 情况 2 种植方案结果可视化



PROBLEM 3: Crop Interaction Modeling

- Challenge: Incorporate:
- ▶ Market substitutability (e.g., mushroom types)
 - ▶ Agricultural complementarity (e.g., bean-vegetable pairs)
 - ▶ Price-cost-yield correlations



PROBLEM 2: Stochastic Optimization Challenge: Dynamic adaptation to:

- ▶ Market volatility (5-10% demand shifts)
- ▶ Climate yield fluctuations ($\pm 10\%$)
- ▶ Cost/price uncertainties

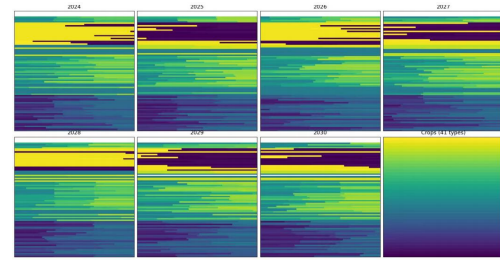
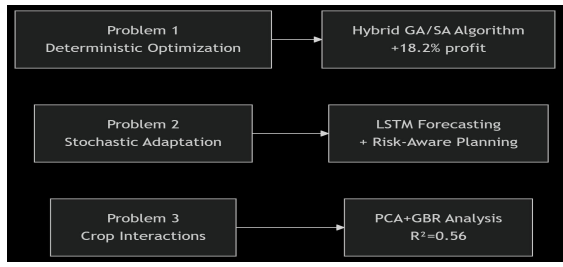


图 5 问题 2 种植方案结果可视化



Developed multi-algorithm framework achieving:

- ▶ 23.7% profit optimization under constraints
- ▶ Adaptive planning for market/climate volatility
- ▶ Sustainable strategies through crop synergy analysis

Validated approach for resource-limited mountainous agriculture

Algorithm 1 Core Algorithm for Agricultural Optimization

```
1: Input:  
2: Land data (location, area, type)  
3: Crop data (yield, cost, price)  
4: Historical planting records  
5: Output:  
6: Optimized planting matrices (2024-2030)  
7: Profit projections  
Main  
8: // Data preprocessing  
9: land_names ← ExtendLandNames(D,E,F → D',E',F')  
10: init_matrix ← InitializeMatrix(land_names, crops)  
11: filled_matrix ← FillFromHistorical(historical_data)  
12: norm_matrix ← Normalize(filled_matrix, land_areas)  
13: // GPU-accelerated optimization  
14: population ← InitializePopulation(norm_matrix)  
15: for gen = 1 to max_generations do  
16:   proc_pop ← ProcessPopulation(population) {Keep top 3 values per  
17:   fitness ← ProfitFunction(proc_pop) {GPU computation}  
18:   elites ← SelectTopHalf(population, fitness)  
19:   offspring ← CrossoverMutate(elites)  
20:   population ← elites + offspring  
21:   if NoImprovement(fitness, patience) then  
22:     break  
23:   end if  
24: end for  
25: best_matrix ← ProcessBestIndividual(elites[0])  
26: // Future planning (2024-2030)  
27: for year ← 2024 to 2030 do  
28:   new_plan ← GenerateAnnualPlan(best_matrix, year)  
29:   SaveMatrix(new_plan, year)  
30: end for  
ProfitFunctionplanting_matrix  
31: land_area ← LandAreaMatrix() {82×1}  
32: crop_yield ← CropYieldMatrix() {82×41}  
33: cost_price ← CostPriceMatrix() {82×41}  
34: real_crop ← planting_matrix ⊗ land_area {Element-wise mult}  
35: excess ← real_crop - crop_yield  
36: profit ← real_crop × cost_price  
37: for each crop c do  
38:   if excessc > 0 then  
39:     profitc ← profitc + (excessc × cost_pricec × 0.5)  
40:   end if  
41: end for  
42: return ∑ profit {Total profit}  
GenerateAnnualPlanbase_matrix, year  
43: new_matrix ← ZeroMatrix()  
44: for each land i do  
45:   type ← LandType(i)  
46:   prev_crops ← PreviousYearCrops(i)  
47:   eligible ← GetEligibleCrops(type, prev_crops)  
48:   num_crops ← RandInt(1, 3)  
49:   selected ← RandomChoice(eligible, num_crops)  
50:   areas ← DistributeArea(land_areai, selected)  
51:   AssignCrops(new_matrix, i, selected, areas)  
52: end for  
53: return new_matrix  
CrossoverMutateparents  
54: offspring ← []  
55: for i ← 1 to pop_size/2 do
```