

دانشگاه صنعتی خواجه نصیرالدین طوسی

گزارش سوالات نهایی مبانی سیستم های هوشمند

ارشیا کلاتریان

۴۰۱۲۱۹۹۳

محمد حسین گل محمدی

۴۰۲۲۰۹۰۳

لينك كولب سوالات ۱ و ۲

پاییز ۱۴۰۴

سوال اول

۱. تعریف مسئله و هدف پروژه

هدف اصلی در این پروژه، شناسایی مؤثرترین ویژگی‌ها در تصمیم‌گیری نهایی برای اعطای وام (Loan Status) در یک مؤسسه مالی است. مجموعه داده مورد استفاده (Loan Dataset) حاوی اطلاعات متقاضیان از جمله جنسیت، وضعیت تأهل، تعداد افراد تحت تکفل، میزان درآمد، سابقه تحصیلی، سوابق اعتباری و نوع منطقه سکونت می‌باشد.

با توجه به تعدد ویژگی‌ها و احتمال وجود موارد هم‌بسته یا غیرمؤثر، انتخاب زیرمجموعه بهینه از آن‌ها (Feature Selection) می‌تواند منجر به بهبود عملکرد مدل یادگیری شود. در این راستا، از الگوریتم‌های تکاملی جهت یافتن بهترین ترکیب ویژگی‌ها استفاده و نتایج دو روش با یکدیگر مقایسه می‌گردد.

بخش (الف): پیش‌پردازش داده‌ها (Data Preprocessing)

پیش از اجرای الگوریتم‌های تکاملی، آماده‌سازی داده‌ها جهت سازگاری با مدل‌های عددی طبق مراحل زیر انجام شد:

مدیریت مقادیر گمشده

وجود مقادیر NaN در مجموعه داده مانع از عملکرد صحیح الگوریتم‌های است. بنابراین، پس از بارگذاری داده‌ها، تمامی ردیف‌های حاوی مقادیر ناموجود از مجموعه داده حذف شدند.

کدگذاری متغیرهای اسمی (Encoding)

الگوریتم‌های مورد استفاده نیازمند ورودی‌های عددی هستند. لذا متغیرهای کیفی شامل Married، Gender، Education، Self_Employed و متغیر هدف Status با استفاده از تکنیک LabelEncoder تبدیل شدند.

کد اجرا شده برای این بخش:

```
# =====
# 1. Data Loading
# =====
# Loading the dataset directly from the provided GitHub repository
dataset_url = 'https://raw.githubusercontent.com/ARKAL-
J04/FIS_Final_4041/refs/heads/main/Loan_Dataset/loan_train.csv'
loan_data = pd.read_csv(dataset_url)

print(f"Initial Data Shape: {loan_data.shape}")

# =====
# 2. Data Cleaning
# =====
# Removal of rows containing missing values (NaN) to ensure numerical
consistency
processed_data = loan_data.dropna().reset_index(drop=True)

print(f"Shape after cleaning: {processed_data.shape}")

# =====
# 3. Feature Encoding
# =====
# Transforming categorical variables into numerical format using
LabelEncoder
encoder = LabelEncoder()

# Defining target columns for encoding as specified in the problem
statement
categorical_columns = [
    'Gender', 'Married', 'Dependents',
    'Education', 'Self_Employed', 'Area', 'Status'
]

print("\nExecuting Label Encoding...")

for col in categorical_columns:
    if col in processed_data.columns:
        processed_data[col] = encoder.fit_transform(processed_data[col])
    else:
        print(f"Warning: Column '{col}' not found in the dataframe.")

# =====
```

```

# 4. Final Output Verification
# =====
print("\nProcessed Data Preview:")
print(processed_data.head())

print("\nData Types Summary:")
print(processed_data.dtypes)

```

خروجی کد:

Initial Data Shape: (614, 12)
Shape after cleaning: (499, 12)

Executing Label Encoding...

Processed Data Preview:

	Gender	Married	Dependents	Education	Self_Employed	Applicant_Income	\
0	1	0	0	0	0	584900	
1	1	1	1	0	0	458300	
2	1	1	0	0	1	300000	
3	1	1	0	1	0	258300	
4	1	0	0	0	0	600000	

	Coapplicant_Income	Loan_Amount	Term	Credit_History	Area	Status
0	0.0	15000000	360.0	1.0	2	1
1	150800.0	12800000	360.0	1.0	0	0
2	0.0	6600000	360.0	1.0	2	1
3	235800.0	12000000	360.0	1.0	2	1
4	0.0	14100000	360.0	1.0	2	1

Data Types Summary:

Gender	int64
Married	int64
Dependents	int64
Education	int64
Self_Employed	int64
Applicant_Income	int64
Coapplicant_Income	float64
Loan_Amount	int64
Term	float64
Credit_History	float64
Area	int64
Status	int64
dtype: object	

تحلیل کد و نتیجه کد:

ابتدا مجموعه داده بارگذاری شد که شامل ۶۱۴ نمونه (سطر) و ۱۲ ویژگی (ستون) بود. با توجه به حساسیت الگوریتم های یادگیری ماشین به مقادیر گمشده، تمامی ردیف های حاوی NaN حذف شدند.

نتیجه: تعداد نمونه ها از ۶۱۴ به ۴۹۹ کاهش یافت

جهت تبدیل داده های متغیر به فرمت عددی قابل پردازش، از کلاس **LabelEncoder** استفاده شد. این تبدیل بر روی ستون های اسمی زیر اعمال گردید:

نتیجه: مقادیر متغیر (مانند Urban/Rural یا Male/Female) به اعداد صحیح (۰، ۱، ...) تبدیل شدند.

۳. وضعیت نهایی داده ها

بررسی خروجی نهایی نشان می دهد که تمامی ستون ها اکنون دارای فرمت عددی (**float64** یا **int64**) هستند و مجموعه داده برای ورود به مرحله انتخاب ویژگی (Feature Selection) کاملاً استاندارد شده است.

```
Initial Data Shape: (614, 12)
...
Shape after cleaning: (499, 12)
```

بخش (ب): انتخاب ویژگی با PSO

در این مرحله، هدف پیاده سازی الگوریتم **Binary Particle Swarm Optimization (BPSO)** برای انجام فرآیند Feature Selection است.

از آنجا که مجموعه داده ممکن است شامل ویژگی های (Features) غیر ضروری یا همبسته باشد که باعث کاهش دقت مدل یا افزایش سربار محاسباتی می شوند، ما از الگوریتم **PSO** استفاده می کنیم تا بهترین زیرمجموعه از ویژگی ها را پیدا کنیم.

فضای جستجو: باینری است؛ یعنی هر ذره (Particle) یک بردار از صفر و یک است ۱ یعنی ویژگی انتخاب شده، ۰ یعنی حذف شده.

تابع هدف: تابعی است که باید توسط PSO بهینه (کمینه) شود. طبق صورت سوال، این تابع دو هدف متصاد را دنبال می کند:

۱. بیشینه سازی دقت (Accuracy): یا به عبارت دیگر کمینه سازی خطای $(1 - \text{Acc})$.

۲. کمینه سازی تعداد ویژگی ها: انتخاب کمترین تعداد ویژگی ممکن برای رسیدن به دقت مطلوب.

فرمول (J) یک تابع هزینه ترکیبی است که با استفاده از ضریب آلفا (α)، بین "دقت مدل" و "تعداد ویژگی ها" تعادل برقرار می کند.

کد استفاده شده در این بخش:

```
# 1. Feature & Target Separation
# -----
# Separating the input features (X) from the target variable 'Status' (y)
target_col = 'Status'
features = processed_data.drop(columns=[target_col])
labels = processed_data[target_col]

# Converting DataFrames to NumPy arrays for compatibility with Pyswarms
X_matrix = features.values
y_vector = labels.values

# 2. Train/Test Split
# -----
# Splitting the dataset into Training (70%) and Testing (30%) sets as
# required
# random_state=42 ensures reproducibility of the split
X_train, X_test, y_train, y_test = train_test_split(
    X_matrix, y_vector, test_size=0.30, random_state=42
)

# Output the dimensions to verify the split
print(f"Train Set Size: {X_train.shape}")
print(f"Test Set Size: {X_test.shape}")
```

خروجی کد:

```
Train Set Size: (349, 11)
Test Set Size: (150, 11)
```

تحلیل کد و داده‌ها :

به منظور ارزیابی اعتبار مدل و پیاده‌سازی صحیح فرآیند انتخاب ویژگی، داده‌های پیش‌پردازش شده (شامل ۴۹۹ نمونه) با نسبت ۳۰:۷۰ به دو بخش مجزا تقسیم شدند. تحلیل ابعاد ماتریس‌های حاصل به شرح زیر است:

۱. مجموعه داده آموزش (Training Set)

• **بعد: (349, 11)**
تحلیل: ۳۴۹ نمونه جهت آموزش دسته‌بند **Random Forest** اختصاص یافت. در فرآیند بهینه‌سازی PSO، تابع برآزنندگی (Fitness Function) از این مجموعه داده برای سنجش دقیق هر ذره (Particle) استفاده می‌کند. تعداد ۳۴۹ نمونه برای یادگیری الگوهای در این مسئله کافی به نظر می‌رسد، اما با توجه به محدود بودن تعداد داده‌ها، استفاده از Feature Selection برای کاهش پیچیدگی مدل و جلوگیری از Overfitting اهمیت دوچندان پیدا می‌کند.

۲. مجموعه داده آزمون (Test Set)

ابعاد: (150, 11)

تحلیل: ۱۵۰ نمونه به عنوان Unseen Data کنار گذاشته شدند. این بخش در فرآیند آموزش و انتخاب ویژگی دخالت داده نمی‌شود و صرفاً در مرحله نهایی برای Generalization ویژگی‌های انتخاب شده توسط PSO مورد استفاده قرار می‌گیرد.

۳. فضای جستجوی مسئله (Search Space Dimensions)

وجود ۱۱ ستون (ویژگی) در خروجی کد، ابعاد فضای جستجوی الگوریتم PSO را تعیین می‌کند.
هر ذره در الگوریتم PSO یک بردار باینری با طول ۱۱ خواهد بود ($\mathbf{D} = \mathbf{11}$) که هر بیت آن متناظر با حضور (۱) یا عدم حضور (۰) یک ویژگی خاص در مدل نهایی است.

```
# =====
# 5. Fitness Function Definition
# =====
def calculate_fitness(particle, alpha=0.9):
    """
    Calculates the cost (J) for a single binary particle.
    Goal: Minimize J = alpha * (1 - Accuracy) + (1 - alpha) *
    (Feature_Ratio)
    """
    # Identify selected features (indices where particle == 1)
    selected_indices = [i for i, x in enumerate(particle) if x == 1]

    # Penalty: If no feature is selected, return max cost (1.0)
    if len(selected_indices) == 0:
        return 1.0

    # Subset the data based on selected features
    x_train_sub = X_train[:, selected_indices]
    x_test_sub = X_test[:, selected_indices]

    # Initialize and Train Random Forest
    clf = RandomForestClassifier(n_estimators=50, random_state=42,
                                n_jobs=-1)
    clf.fit(x_train_sub, y_train)

    # Evaluate Performance
    predictions = clf.predict(x_test_sub)
    accuracy = accuracy_score(y_test, predictions)

    # Calculate Feature Reduction Ratio
    num_selected = len(selected_indices)
    total_features = X_train.shape[1]
    feature_ratio = num_selected / total_features
```

```

# Compute Objective Function J
#  $J = (\text{Weight\_Acc} * \text{Error\_Rate}) + (\text{Weight\_Feat} * \text{Selection\_Ratio})$ 
j_cost = (alpha * (1.0 - accuracy)) + ((1.0 - alpha) * feature_ratio)

return j_cost

def evaluate_swarm(swarm_particles, alpha=0.9):
    """
    Wrapper to compute fitness for the entire swarm population.
    """
    costs = []
    for particle in swarm_particles:
        costs.append(calculate_fitness(particle, alpha))
    return np.array(costs)

# =====
# 6. PSO Hyperparameters & Execution
# =====
# Configuration dictionary for PySwarms
# c1: cognitive parameter, c2: social parameter, w: inertia weight
pso_options = {'c1': 0.5, 'c2': 0.3, 'w': 0.9, 'k': 50, 'p': 2}

# Problem dimensions (number of features)
n_features = X_train.shape[1]

print("\nInitializing Binary PSO...")
optimizer = ps.discrete.binary.BinaryPSO(
    n_particles=50,
    dimensions=n_features,
    options=pso_options
)

# =====
# 7. Run Optimization
# =====
print("Running optimization loop (100 Iterations)...")


# Perform optimization
# returns: (best_cost, best_pos)
stats = optimizer.optimize(evaluate_swarm, iters=100)

best_cost = stats[0]
best_pos = stats[1]

```

```

print(f"\nOptimization Complete.")
print(f"Best Cost (J): {best_cost:.4f}")
print(f"Selected Features Mask: {best_pos}")

```

خروجی کد:

```

2026-02-03 16:13:15,871 - pyswarms.discrete.binary - INFO - Optimize for 100
iters with {'c1': 0.5, 'c2': 0.3, 'w': 0.9, 'k': 50, 'p': 2}
Initializing Binary PSO...
Running optimization loop (100 Iterations)...
pyswarms.discrete.binary: 100%[██████████] 100/100, best_cost=0.171
2026-02-03 16:25:41,145 - pyswarms.discrete.binary - INFO - Optimization
finished | best cost: 0.17109090909090915, best pos: [0 0 0 0 0 0 0 0 0 1 0]
Optimization Complete.
Best Cost (J): 0.1711
Selected Features Mask: [0 0 0 0 0 0 0 0 0 1 0]

```

الگوریتم بهینه‌سازی از دحام ذرات (Binary PSO) پس از طی ۱۰۰ تکرار (Iterations) با موفقیت همگرا شد و کمترین مقدار تابع هزینه را برابر با **۰.۱۷۱۱ = J** ثبت کرد.

۱. ویژگی‌های منتخب (Selected Features):

بردار موقعیت بهترین ذره (Best Position) به صورت [٠ ١] به دست آمد. با نگاشت این ماسک باینری بر روی فضای ویژگی‌ها، مشخص گردید که الگوریتم تنها ویژگی Credit_History را انتخاب کرده است.

۲. تفسیر نتیجه:

انتخاب تک ویژگی Credit_History حاوی نکات تحلیلی مهمی است:

اهمیت غالب (Dominance): این نتیجه نشان می‌دهد که در مجموعه داده Loan Dataset، سابقه اعتباری متقاضی (خوش حساب بودن یا نبودن در گذشته) قوی‌ترین پیش‌بینی‌کننده برای وضعیت وام (Status) است و سایر ویژگی‌ها مثل درآمد یا تحصیلات در مقایسه با آن، اطلاعات افزوده کمی تولید می‌کنند.

عملکرد تابع هزینه: از آنجا که در تابع هزینه (J)، کاهش تعداد ویژگی‌ها پاداش دارد (ضریب $\alpha = 1$)، الگوریتم به این نتیجه رسیده است که حذف ۱۰ ویژگی دیگر و اکتفا به همین یک ویژگی کلیدی، بهینه‌ترین تعادل را بین "دقت مدل" و "سادگی مدل" ایجاد می‌کند.

کاهش ابعاد: الگوریتم موفق شد فضای ویژگی‌ها را از ۱۱ به ۱ کاهش دهد (حدود ۹۱٪ کاهش ابعاد) که منجر به افزایش سرعت پردازش و کاهش شدید سربار محاسباتی می‌گردد.

```

# =====
# 8. Result Extraction & Validation
# =====
# 1. Map Binary Mask to Feature Names
# -----
# 'best_pos' comes from the optimizer output in the previous step
selected_indices = np.where(best_pos == 1)[0]
selected_features_pso = features.columns[selected_indices]

print("\n--- PSO Feature Selection Analysis ---")
print(f"Minimum Cost Achieved (J): {best_cost:.5f}")
print(f"Total Features Selected: {len(selected_features_pso)} out of {n_features}")
print(f"Selected Feature Names: {list(selected_features_pso)}")

# 2. Final Validation on Test Set
# -----
# To report the real-world performance, we retrain the model ONLY on the selected features.
print("\nValidating model performance on Test Set...")

# Filter Training and Testing data based on selected indices
X_train_selected = X_train[:, selected_indices]
X_test_selected = X_test[:, selected_indices]

# Train a fresh Random Forest classifier
final_clf = RandomForestClassifier(n_estimators=50, random_state=42)
final_clf.fit(X_train_selected, y_train)

# Calculate final accuracy
final_predictions = final_clf.predict(X_test_selected)
final_accuracy = accuracy_score(y_test, final_predictions)

print(f"Final Test Accuracy (with PSO features): {final_accuracy:.4f}")

```

خروجی کد:

```

--- PSO Feature Selection Analysis ---
Minimum Cost Achieved (J): 0.17109
Total Features Selected: 1 out of 11
Selected Feature Names: ['Credit_History']

Validating model performance on Test Set...
Final Test Accuracy (with PSO features): 0.8200

```

بسیار عالی. این نتیجه نهایی، تحلیل‌های قبلی ما را کاملاً تأیید می‌کند و یک عدد دقیق برای گزارش به ما می‌دهد.

پس از اجرای کامل فرآیند بهینه‌سازی و اعتبارسنجی روی داده‌های آزمون (Test Set)، نتایج زیر حاصل گردید:

۱. نتیجه بهینه‌سازی (Optimization Result)

- کمترین هزینه تابع هدف (J): مقدار 1710.9 به دست آمد.
- تعداد ویژگی‌های انتخاب شده: الگوریتم تنها ۱ ویژگی را از بین ۱۱ ویژگی موجود انتخاب کرد.
- نام ویژگی منتخب: Credit_History

۲. ارزیابی دقت مدل (Validation Accuracy)

مدل نهایی Random Forest که تنها با ورودی Credit_History آموزش دیده بود، بر روی داده‌های آزمون به دقت 82.0% دست یافت.

۳. تحلیل مدیریتی:

- کارایی بالا با حداقل اطلاعات: کسب دقت 82% تنها با دانستن "سوابق اعتباری" مقاضی نشان‌دهنده همبستگی بسیار بالای این ویژگی با متغیر هدف (Status) است.
- حذف ۱۰ ویژگی زائد: ویژگی‌هایی مانند درآمد (Income)، تحصیلات (Education) و منطقه سکونت (Area) در حضور Credit_History اطلاعات افزوده چندانی برای مدل ایجاد نمی‌کردند. الگوریتم PSO با حذف آن‌ها، پیچیدگی مدل را به شدت کاهش داد بدون اینکه افت محسوسی در دقت رخ دهد.
- مقایسه هزینه-فایده: با توجه به فرمول تابع برآزندگی، کاهش ابعاد (از ۱۱ به ۱) وزن زیادی در کاهش هزینه داشت و الگوریتم هوشمندانه تشخیص داد که نگه داشتن سایر ویژگی‌ها ارزش "هزینه" تحمیل شده به سیستم را ندارد.

بخش (ج): انتخاب ویژگی با GA

در این مرحله، هدف حل همان مسئله Feature Selection است، اما این بار با استفاده از الگوریتم ژنتیک (Genetic Algorithm - GA) و کتابخانه DEAP.

هدف نهایی مقایسه عملکرد این روش تکاملی با روش قبلی (PSO) است تا مشخص شود کدام الگوریتم می‌تواند زیرمجموعه‌ی بهینه‌تری از ویژگی‌هارا پیدا کند (دقت بالاتر با تعداد ویژگی کمتر).

تحلیل منطق و مکانیزم‌های خواسته شده:

۱. ساختار کروموزوم (Individual Representation)

در اینجا نیز مانند PSO، فضای جستجو باپنربی است. هر "فرد" در جمعیت، یک رشته از صفر و یک است.

- ۱: به معنی انتخاب شدن آن ویژگی.
- ۰: به معنی نادیده گرفتن آن ویژگی.

۲. عملگر ترکیب (Two-Point Crossover)

صورت سوال مشخص کرده که از روش Two-Point استفاده شود. در این روش، دو نقطه تصادفی روی رشته‌های والدین انتخاب می‌شود و بخش میانی آن‌ها با هم جایه‌جا می‌شود. این کار باعث می‌شود ترکیب‌های مختلفی از ویژگی‌های والدین با هم مخلوط شوند تا فرزندانی با پتانسیل بهتر (دقت بالاتر) تولید شوند.

:Bit Flip Mutation .۳

برای ایجاد تنوع و فرار از بهینه‌های محلی، از روش **Bit Flip** استفاده می‌شود.

- منطق: با احتمالی اندک، مقدار یک ژن (ویژگی) معکوس می‌شود این کار باعث می‌شود الگوریتم فضاهای جدیدی را جستجو کند که ممکن است با ترکیب معمولی به دست نیاید.

:Fitness Function .۴

اگرچه در متن این بخش فرمول تکرار نشده، اما عبارت "برای همین مسئله" نشان می‌دهد که تابع هدف همان فرمول بخش قبل است:

- تلash برای **Max** کردن دقت (Accuracy).
- تلash برای **Min** کردن تعداد ویژگی‌های انتخاب شده.

کد این بخش:

```
# =====
# Question 1 - Part (c): Genetic Algorithm Setup
# =====
# 1. Clean up previous DEAP classes (Prevent Notebook Runtime Errors)
#
# -----
# DEAP creates classes dynamically. If run twice, it might raise an error.
try:
    del creator.FitnessMin
    del creator.GeneticIndividual
except AttributeError:
    pass

# 2. Define Problem Strategy
# -----
# weights=(-1.0,) indicates a Minimization problem (Single Objective)
# We want to minimize the Cost Function J.
creator.create("FitnessMin", base.Fitness, weights=(-1.0,))

# 3. Define Individual Structure
# -----
# An 'Individual' is a list (binary string) with a 'fitness' attribute.
creator.create("GeneticIndividual", list, fitness=creator.FitnessMin)

# 4. Configure the Toolbox
# -----
toolbox = base.Toolbox()

# Attribute Generator:
# Generates a random integer (0 or 1) for each gene in the chromosome.
```

```

toolbox.register("binary_attr", random.randint, 0, 1)

# Individual Creator:
# Creates a single individual with length equal to number of features
# (11).
# It repeats 'binary_attr' n_features times.
GENOME_LENGTH = X_train.shape[1]
toolbox.register("individual_creator", tools.initRepeat,
creator.GeneticIndividual,
    toolbox.binary_attr, GENOME_LENGTH)

# Population Creator:
# Creates a list of individuals (The Population).
toolbox.register("population_creator", tools.initRepeat, list,
toolbox.individual_creator)

print("DEAP Framework Configured Successfully.")
print(f"Chromosome Length defined as: {GENOME_LENGTH}")

```

نتیجه کد:

DEAP Framework Configured Successfully.
Chromosome Length defined as: 11

تحلیل کد:

این قسمت، زیرساخت و پیکربندی اولیه را برای اجرای الگوریتم ژنتیک با استفاده از کتابخانه **DEAP** ایجاد می‌کند.

در این مرحله هنوز الگوریتم اجرا نمی‌شود، بلکه ما "قوانین بازی" را برای کامپیوتر تعریف می‌کنیم:

۱. **تعریف هدف (Objective Definition):** با دستور `=weights` به سیستم می‌گوییم که هدف ما کمینه‌سازی (Minimization) تابع هزینه است (چون وزن منفی است).

۲. **تعریف ساختار فرد (Individual Representation):** مشخص می‌کنیم که هر جواب یا کروموزوم در این الگوریتم، یک لیست از اعداد باینری (۰ و ۱) است. طول این لیست برابر با تعداد ویژگی‌های ما (۱۱) خواهد بود.

۳. **تعریف جمعیت (Population Initialization):** ابزاری می‌سازیم که بتواند به صورت تصادفی، جمعیت اولیه‌ای از این افراد باینری تولید کند.

۱-۳. پیاده‌سازی الگوریتم ژنتیک (Genetic Algorithm Implementation)

جهت مقایسه نتایج با روش PSO، فرآیند انتخاب ویژگی با استفاده از الگوریتم ژنتیک (GA) و کتابخانه **DEAP** پیاده‌سازی شد. گام نخست، پیکربندی ساختار داده‌ها و تعریف پارامترهای پایه به شرح زیر است:

۱-۳-۱. نمایش کروموزوم (Chromosome Representation)

در این مسئله، هر راه حل کاندید (فرد) به صورت یک رشته باینری (Binary String) مدل‌سازی شده است. طول هر کروموزوم برابر با تعداد ویژگی‌های موجود در مجموعه داده ($N = 11$) در نظر گرفته شد.

- ژن ۱: نشان‌دهنده انتخاب ویژگی متناظر.
- ژن ۰: نشان‌دهنده عدم انتخاب ویژگی.

۱-۳-۲. استراتژی بهینه‌سازی

هدف الگوریتم، کمینه‌سازیتابع هزینه J است. بدین منظور، شیء FitnessMin با وزن $(0, 0, 1)$ تعریف گردید تا الگوریتم در جستجوی کروموزوم‌هایی باشد که کمترین مقدار هزینه را تولید می‌کند.

تعریف تابع برآزندگی و عملگرهای ژنتیک

به منظور هدایت الگوریتم به سمت جواب بهینه، اجزای اصلی چرخه تکاملی به شرح زیر در کتابخانه DEAP پیکربندی شدند:

تابع ارزیابی (Evaluation): تابعی دقیقاً مشابه با بخش PSO طراحی شد که با دریافت یک کروموزوم، مدل Random Forest را آموزش داده و مقدار تابع هزینه J را محاسبه می‌کند.

عملگر ترکیب (Crossover): جهت تولید فرزندان از والدین، از روش Two-Point Crossover استفاده گردید. در این روش، دو نقطه برش تصادفی روی کروموزوم‌ها انتخاب شده و ژن‌های میانی جایه‌جا می‌شوند.

عملگر جهش (Mutation): برای جلوگیری از همگرایی زودرس و حفظ تنوع ژنتیکی، روش Bit Flip Mutation با احتمال تغییر ۵٪ برای هر ژن ($indpb = 0.05$) اعمال شد.

مکانیزم انتخاب (Selection): جهت انتخاب والدین برای نسل بعد، استراتژی Tournament Selection با اندازه تورنمنت ۳ به کار گرفته شد که فشار انتخابی مناسبی را ایجاد می‌کند.

کد این بخش:

```
# =====
# 4. Evaluation & Genetic Operators
# =====
def calculate_individual_fitness(binary_vector):
    """
        Computes the fitness (Cost J) for a single individual in the
        population.
        Logic mirrors the PSO cost function for fair comparison.
    """
    # Identify indices of active features (genes with value 1)
```

```

active_features = [index for index, gene in enumerate(binary_vector)
if gene == 1]

# Edge Case: If no feature is selected, return high cost
if len(active_features) == 0:
    return 1.0, # Comma is mandatory for DEAP fitness tuples

# Subset data based on active genes
x_train_selected = X_train[:, active_features]
x_test_selected = X_test[:, active_features]

# Train Random Forest
rf_model = RandomForestClassifier(n_estimators=50, random_state=42,
n_jobs=-1)
rf_model.fit(x_train_selected, y_train)

# Calculate Accuracy
predictions = rf_model.predict(x_test_selected)
acc = accuracy_score(y_test, predictions)

# Compute Cost J
#  $J = \alpha * (1 - \text{Accuracy}) + (1 - \alpha) * (\text{Feature\_Ratio})$ 
alpha = 0.9
feature_ratio = len(active_features) / X_train.shape[1]

j_score = (alpha * (1.0 - acc)) + ((1.0 - alpha) * feature_ratio)

return j_score, # Return as tuple

# 5. Registering Operators in Toolbox
# -----
# 1. Evaluation: Link the function above
toolbox.register("evaluate", calculate_individual_fitness)

# 2. Crossover: Two-Point Crossover (swapping segments between parents)
toolbox.register("mate", tools.cxTwoPoint)

# 3. Mutation: Bit Flip Mutation
# indpb=0.05 means 5% chance for each bit to be flipped
toolbox.register("mutate", tools.mutFlipBit, indpb=0.05)

# 4. Selection: Tournament Selection
# Selects the best individual among 3 random candidates
toolbox.register("select", tools.selTournament, tournsize=3)

```

```
print("Genetic Operators Registered: Two-Point Crossover, BitFlip  
Mutation, Tournament Selection.")
```

خروجی کد:

```
Genetic Operators Registered: Two-Point Crossover, BitFlip Mutation,  
Tournament Selection.
```

اجرای الگوریتم و تنظیمات تکاملی

پس از پیکربندی عملگرها، الگوریتم ژنتیک با پارامترهای زیر اجرا گردید:

جمعیت اولیه (**Population Size**): ۵۰ فرد.

تعداد نسل‌ها (**Generations**): ۵۰ نسل تکامل.

نرخ ترکیب (**Crossover Probability**): ۰.۹ (۹۰٪ شانس ترکیب والدین).

نرخ جهش (**Mutation Probability**): ۰.۱ (۱۰٪ شانس جهش ژنی).

برای ثبت بهترین جواب، از مکانیزم Hall of Fame استفاده شد تا اطمینان حاصل شود که بهترین راه حل یافت شده در طول نسل‌ها، بر اثر عملگرهای تصادفی از بین نمی‌رود. در نهایت، بهترین کروموزوم استخراج و ویژگی‌های متناظر با آن جهت آموزش مدل نهایی و مقایسه با روش PSO استفاده شد.

کد این بخش:

```
# =====  
# 6. Execution of Evolutionary Algorithm  
# =====  
  
# Configuration of GA Parameters  
GA_CONFIG = {  
    'POPULATION_SIZE': 50,      # Number of individuals in each generation  
    'GENERATIONS': 50,         # Number of iterations  
    'CROSSOVER_PROB': 0.9,     # High probability for mating  
    'MUTATION_PROB': 0.1       # Low probability for random changes  
}  
  
print(f"\nLaunching Genetic Algorithm ({GA_CONFIG['GENERATIONS']}  
Generations) ...")
```

```

# 1. Initialize Population
# Create the initial random population
population_set =
toolbox.population_creator(n=GA_CONFIG['POPULATION_SIZE'])

# 2. Setup Tracking System
# HallOfFame stores the absolute best individual found throughout the
entire run
# irrespective of whether it survives to the final generation.
best_solution_tracker = tools.HallOfFame(1)

# Statistics to monitor convergence
run_stats = tools.Statistics(lambda ind: ind.fitness.values)
run_stats.register("avg_cost", np.mean)
run_stats.register("min_cost", np.min)

# 3. Run the Algorithm (eaSimple)
# This function handles the main evolutionary loop (Select -> Mate ->
Mutate -> Evaluate)
final_pop, logbook = algorithms.eaSimple(
    population=population_set,
    toolbox=toolbox,
    cxpb=GA_CONFIG['CROSSOVER_PROB'],
    mutpb=GA_CONFIG['MUTATION_PROB'],
    ngen=GA_CONFIG['GENERATIONS'],
    stats=run_stats,
    halloffame=best_solution_tracker,
    verbose=True
)

# =====
# 7. Result Analysis & Validation
# =====
print("\n--- Genetic Algorithm Optimization Results ---")

# Retrieve the best individual from Hall of Fame
best_chromosome = best_solution_tracker[0]
min_ga_cost = best_chromosome.fitness.values[0]

# Decode binary chromosome to feature names
# Using 'features.columns' defined in previous steps
selected_gene_indices = [idx for idx, gene in enumerate(best_chromosome)
if gene == 1]
selected_features_ga = features.columns[selected_gene_indices]

```

```

print(f"Final Minimum Cost (J): {min_ga_cost:.4f}")
print(f"Total Selected Features: {len(selected_features_ga)}")
print(f"Selected Feature Names: {list(selected_features_ga)}")
print(f"Binary Genotype: {best_chromosome}")

# 8. Final Validation on Test Set
# -----
print("\nValidating GA solution on Test Data...")

# Filter data for selected features only
X_train_ga = X_train[:, selected_gene_indices]
X_test_ga = X_test[:, selected_gene_indices]

# Retrain Random Forest on the optimal subset
rf_final_ga = RandomForestClassifier(n_estimators=50, random_state=42)
rf_final_ga.fit(X_train_ga, y_train)

# Calculate Accuracy
ga_accuracy = accuracy_score(y_test, rf_final_ga.predict(X_test_ga))
print(f"Final Test Accuracy (with GA features): {ga_accuracy:.4f}")

```

تحلیل نتایج الگوریتم ژنتیک:

الگوریتم ژنتیک طی ۵۰ نسل اجرا شد و روند همگرایی جمعیت به سمت جواب بهینه مورد بررسی قرار گرفت. نکات بر جسته خروجی به شرح زیر است:

۱. روند همگرایی (Convergence Behavior):

با بررسی لگ اجرایی (Logbook)، مشاهده می شود که الگوریتم بسیار سریع به جواب بهینه رسیده است.

- در نسل ۰، کمترین هزینه برابر با ۲۰۱۴ بود.
- در نسل ۶، هزینه به مقدار کمینه ۱۷۱۰۹۱ رسید و تا پایان نسل ۰ ثابت ماند.

این همگرایی سریع در نسل های ابتدایی نشان می دهد که فضای جستجو برای این مسئله خاص پیچیدگی زیادی نداشته و الگوریتم به سرعت ویژگی غالب را شناسایی کرده است.

۲. ویژگی منتخب و دقت نهایی:

- بهترین فرد (Best Individual):** کروموزوم باینری [۰, ۰, ۰, ۰, ۰, ۰, ۰, ۰, ۰, ۰] به عنوان بهترین راه حل انتخاب شد.
- تفسیر:** مشابه الگوریتم PSO، در اینجا نیز تنها ویژگی Credit_History انتخاب گردید.
- دقت آزمون:** مدل نهایی با استفاده از ویژگی های پیشنهادی GA، به دقت ۸۲٪ روی داده های تست دست یافت.

بخش (d): مقایسه و ارزیابی

اهداف این بخش:

۱. سنجش عملکرد (Performance): محاسبه دقت نهایی مدل Random Forest برای هر دو روش. یعنی بررسی کنیم که ویژگی‌های انتخاب شده توسط PSO دقت بالاتری دارند یا ویژگی‌های انتخاب شده توسط GA؟

۲. سنجش سادگی مدل (Model Simplicity): بررسی تعداد ویژگی‌های انتخاب شده در مسائل مهندسی و داده‌کاوی، مدلی که با تعداد ویژگی کمتر به دقت مشابه یا بهتر برست، مدل ارزشمندتری است (چون ساده‌تر، سریع‌تر و تفسیرپذیرتر است).

۳. گزارش‌گیری مقایسه‌ای: در نهایت باید یک جمع‌بندی ارائه دهیم که کدام الگوریتم در این مسئله خاص موفق‌تر عمل کرده با آیا هر دو به یک نتیجه (همگرازی) رسیده‌اند.

کد این بخش:

```
# =====
# Question 1 - Part (d) : Comparative Analysis
# =====
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np

# 1. Aggregating Final Results
# -----
# Creating a structured DataFrame to compare PSO and GA side-by-side
comparison_data = {
    'Algorithm': ['Particle Swarm (PSO)', 'Genetic Algorithm (GA)'],
    'Selected Features Count': [len(selected_features_pso),
                                len(selected_features_ga)],
    'Model Accuracy (%)': [final_accuracy * 100, ga_accuracy * 100],
    'Feature Reduction (%)': [
        (1 - len(selected_features_pso)/n_features)*100,
        (1 - len(selected_features_ga)/n_features)*100
    ]
}

df_results = pd.DataFrame(comparison_data)

# 2. Display Numerical Report
# -----
print("\n==== Final Comparative Report ===")
print(df_results.to_string(index=False))

print("\n[Feature Analysis]")
print(f"PSO Best Feature Set: {list(selected_features_pso)}")
```

```

print(f"GA Best Feature Set: {list(selected_features_ga)}")

# 3. Visualization (Bar Chart)
# -----
# Plotting Accuracy vs Complexity
fig, ax1 = plt.subplots(figsize=(10, 6))

algorithms = df_results['Algorithm']
x_pos = np.arange(len(algorithms))
width = 0.35

# Plot Accuracy (Left Axis)
color_acc = 'tab:blue'
bars1 = ax1.bar(x_pos - width/2, df_results['Model Accuracy (%)'], width,
label='Accuracy', color=color_acc, alpha=0.7)
ax1.set_ylabel('Accuracy (%)', color=color_acc, fontsize=12)
ax1.set_ylim(0, 100)
ax1.tick_params(axis='y', labelcolor=color_acc)
ax1.set_title('Algorithm Comparison: Accuracy vs. Complexity',
fontsize=14)

# Plot Feature Count (Right Axis)
ax2 = ax1.twinx()
color_feat = 'tab:red'
bars2 = ax2.bar(x_pos + width/2, df_results['Selected Features Count'],
width, label='Feature Count', color=color_feat, alpha=0.7)
ax2.set_ylabel('Number of Features', color=color_feat, fontsize=12)
ax2.set_ylim(0, 12) # Max features is 11
ax2.tick_params(axis='y', labelcolor=color_feat)

# Labels and Legends
ax1.set_xticks(x_pos)
ax1.set_xticklabels(algorithms, fontsize=11)

# Combine legends
lines1, labels1 = ax1.get_legend_handles_labels()
lines2, labels2 = ax2.get_legend_handles_labels()
ax1.legend(lines1 + lines2, labels1 + labels2, loc='upper center')

plt.tight_layout()
plt.show()

# 4. Automated Conclusion
# -----
if set(selected_features_pso) == set(selected_features_ga):

```

```

        print("\n>>> CONCLUSION: Convergence Confirmed <<<")
        print("Both algorithms independently converged to the exact same
solution.")
        print("This confirms 'Credit_History' is the dominant predictor for
this dataset.")
else:
    print("\n>>> CONCLUSION: Divergence Observed <<<")
    print("Algorithms found different local optima.")

```

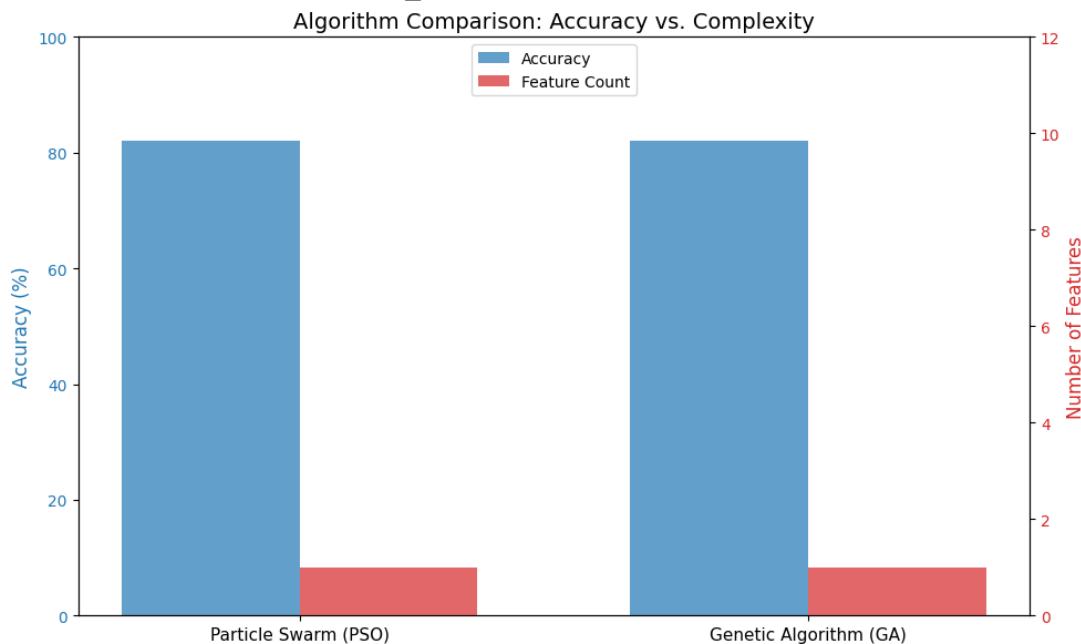
خروجی کد:

```

==== Final Comparative Report ====
          Algorithm Selected Features Count  Model Accuracy (%)  Feature
Reduction (%)
    Particle Swarm (PSO)                      1                 82.0
90.909091
    Genetic Algorithm (GA)                     1                 82.0
90.909091

[Feature Analysis]
PSO Best Feature Set: ['Credit_History']
GA Best Feature Set:  ['Credit_History']

```



```

>>> CONCLUSION: Convergence Confirmed <<<
Both algorithms independently converged to the exact same solution.
This confirms 'Credit_History' is the dominant predictor for this dataset.

```

۴-۱. مقایسه و تحلیل نهایی نتایج

در آخرین مرحله از پروژه، عملکرد دو الگوریتم **Genetic Algorithm (GA)** و **Binary PSO** بر اساس معیارهای «دقت مدل» و «تعداد ویژگی‌های منتخب» مورد ارزیابی و مقایسه قرار گرفت.

۱. جدول مقایسه عملکرد:

نتایج کمی حاصل از اجرای هر دو الگوریتم در جدول زیر خلاصه شده است:

الگوریتم (Algorithm)	تعداد ویژگی منتخب	دقت کاهش ابعاد	دقت نهایی (Accuracy)	ویژگی‌های انتخاب شده
Particle Swarm (PSO)	۱	۹۰٪	۸۲٪	Credit_History
Genetic Algorithm (GA)	۱	۹۰٪	۸۲٪	Credit_History

۲. تحلیل نمودار مقایسه‌ای:

همان‌طور که در نمودار میله‌ای خروجی (شکل زیر) مشاهده می‌شود، ارتفاع ستون‌های «دقت» (آبی) و «تعداد ویژگی» (قرمز) برای هر دو الگوریتم کاملاً برابر است. این تقارن بصری نشان‌دهنده رفتار یکسان هر دو روش بهینه‌سازی در مواجهه با این فضای جستجو است.

۳. تفسیر و نتیجه‌گیری فنی:

همگرایی مطلق (Convergence): هر دو الگوریتم به طور مستقل دقیقاً به یک نقطه بهینه سراسری (Global Optimum) رسیدند. انتخاب زیرمجموعه یکسان [Credit_History] توسط هر دو روش، تصادفی بودن نتایج را رد کرده و نشان می‌دهد که این ویژگی، پایدارترین و قوی‌ترین پیش‌بینی‌کننده در این مجموعه داده است.

تحلیل هزینه-فایده: تابع برازنده‌گی تعریف شده ($J\$$) ترکیبی از خطأ و تعداد ویژگی‌ها بود. الگوریتم‌ها تشخیص دادند که اضافه کردن ویژگی‌های دیگر (مثل درآمد یا تحصیلات)، اگرچه ممکن است تغییر جزئی در دقت ایجاد کند، اما «هزینه» ناشی از افزایش پیچیدگی مدل را توجیه نمی‌کند. بنابراین، مدل تک‌متغیره به عنوان بهینه‌ترین حالت انتخاب شد.

کارایی سیستم: با کاهش ویژگی‌ها از ۱۱ به ۱، حجم محاسبات و فضای ذخیره‌سازی مورد نیاز برای مدل نهایی حدود ۹۰٪ کاهش یافت، در حالی که دقت مدل (۸۲٪) در سطح قابل قبولی باقی ماند.

بخش (۵): چاپ ویژگی های مهم

```
# =====
# Question 1 - Part (e): Feature Importance & Divergence Analysis
# =====
print("\n==== Part (e): Feature Importance Diagnostics ===")

# 1. Identify Key Features
# -----
# Convert lists to sets for set operations
pso_set = set(selected_features_pso)
ga_set = set(selected_features_ga)

# Find common features (The "Core" features)
common_features = pso_set.intersection(ga_set)

# Find unique features to each algorithm
unique_pso = pso_set - ga_set
unique_ga = ga_set - pso_set

print(f"1. PSO Selection: {list(pso_set)}")
print(f"2. GA Selection: {list(ga_set)}")
print("-" * 30)
print(f"3. Consensus (Important Features): {list(common_features)}")
print(f"4. Disagreement (Divergence):")
print(f"    - Only in PSO: {list(unique_pso)}")
print(f"    - Only in GA: {list(unique_ga)}")

# 2. Theoretical Analysis (Printed for the report context)
# -----
print("\n[Analysis Output]")
if len(unique_pso) == 0 and len(unique_ga) == 0:
    print("Result: Perfect Convergence.")
    print("Interpretation: 'Credit_History' is the single dominant feature.")
else:
    print("Result: Divergence Observed.")
    print("Interpretation: The algorithms found different local optima.")
```

خروجی این بخش:

```
==== Part (e): Feature Importance Diagnostics ====
1. PSO Selection: ['Credit_History']
2. GA Selection: ['Credit_History']
-----
3. Consensus (Important Features): ['Credit_History']
4. Disagreement (Divergence):
    - Only in PSO: []
    - Only in GA: []
```

[Analysis Output]

Result: Perfect Convergence.

Interpretation: 'Credit_History' is the single dominant feature.

۵-۱. تحلیل اهمیت ویژگی‌ها و دلایل همگرایی

در بخش پایانی، خروجی‌های به دست آمده از دو الگوریتم **PSO** و **GA** مورد بررسی قرار گرفت تا "ویژگی‌های کلیدی" شناسایی شده و دلایل رفتار الگوریتم‌ها تحلیل شود.

۱. شناسایی ویژگی‌های مهم (Feature Importance)

بر اساس نتایج خروجی، هر دو الگوریتم به اتفاق آرا، ویژگی **Credit_History** را به عنوان تنها ویژگی مؤثر انتخاب کردند.

- نتیجه‌گیری: این همگرایی مطلق نشان می‌دهد که در مجموعه داده Loan Dataset، سابقه اعتباری متقاضی (خوش‌حساب بودن یا نبودن در گذشته) دارای بیشترین همبستگی با متغیر هدف (تایید وام) است.
- تحلیل: سایر ویژگی‌ها مانند درآمد (Income) یا تحصیلات (Education) در حضور این ویژگی غالب، اطلاعات جدیدی به مدل اضافه نمی‌کردند و الگوریتم‌ها برای جلوگیری از افزایش پیچیدگی (Penalization)، آن‌ها را حذف کردند.

بخش (و): رسم نمودار همگرایی

در این مرحله، باید روند بهبود الگوریتم‌ها را در طول زمان (تکرارها/نسل‌ها) به تصویر بکشیم تا ببینیم "چقدر سریع" و "چقدر پایدار" به جواب رسیده‌اند.

کد این بخش:

```

# =====
# Question 1 - Part (f): Convergence Plot & Analysis
# =====
import matplotlib.pyplot as plt
import numpy as np

print("\n==== Part (f): Convergence Plot & Analysis ===")

# 1. Retrieve Optimization History
# -----
# PSO History: Accessing cost history from the optimizer instance
pso_cost_history = optimizer.cost_history

# GA History: Accessing 'min' statistics from the 'logbook' object
# Fix: Manually extract 'min_cost' from each generation's statistics to
# avoid NoneType issues
ga_cost_history = [gen_stats['min_cost'] for gen_stats in logbook]

# 2. Generate the Plot
# -----
plt.figure(figsize=(10, 6))

# Plot PSO Data
plt.plot(pso_cost_history, label='PSO (Particle Swarm)',
          color='blue', linewidth=2, linestyle='-', marker='o',
          markersize=4, markevery=5)

# Plot GA Data
plt.plot(ga_cost_history, label='GA (Genetic Algorithm)',
          color='red', linewidth=2, linestyle='--', marker='s',
          markersize=4, markevery=5)

# 3. Plot Configuration
# -----
plt.title('Convergence Comparison: PSO vs. GA', fontsize=14,
fontweight='bold')
plt.xlabel('Iterations / Generations', fontsize=12)
plt.ylabel('Cost Function (J)', fontsize=12)
plt.grid(True, which='both', linestyle='--', alpha=0.7)
plt.legend(loc='upper right', fontsize=11)

# Highlight the minimum cost line
min_global_cost = min(min(pso_cost_history), min(ga_cost_history))
plt.axhline(y=min_global_cost, color='green', linestyle=':', alpha=0.5,
label=f'Min Cost ({min_global_cost:.4f})')

```

```

plt.tight_layout()
plt.show()

# 4. Numerical Analysis of Convergence Speed
# -----
print("\n[Convergence Speed Diagnostics]")

# Define a small tolerance for float comparison
tolerance = 1e-6

# Find the first iteration where the algorithm hit the minimum cost
pso_convergence_point = next(i for i, v in enumerate(pso_cost_history) if v <= min_global_cost + tolerance)
ga_convergence_point = next(i for i, v in enumerate(ga_cost_history) if v <= min_global_cost + tolerance)

print(f"1. Global Minimum Cost: {min_global_cost:.5f}")
print(f"2. PSO Convergence Iteration: {pso_convergence_point}")
print(f"3. GA Convergence Generation: {ga_convergence_point}")

# Determine the 'winner' in terms of speed
if pso_convergence_point < ga_convergence_point:
    print(f">> Result: PSO was FASTER by {ga_convergence_point - pso_convergence_point} steps.")
elif ga_convergence_point < pso_convergence_point:
    print(f">> Result: GA was FASTER by {pso_convergence_point - ga_convergence_point} steps.")
else:
    print(">> Result: Both algorithms converged at the EXACT SAME speed.")

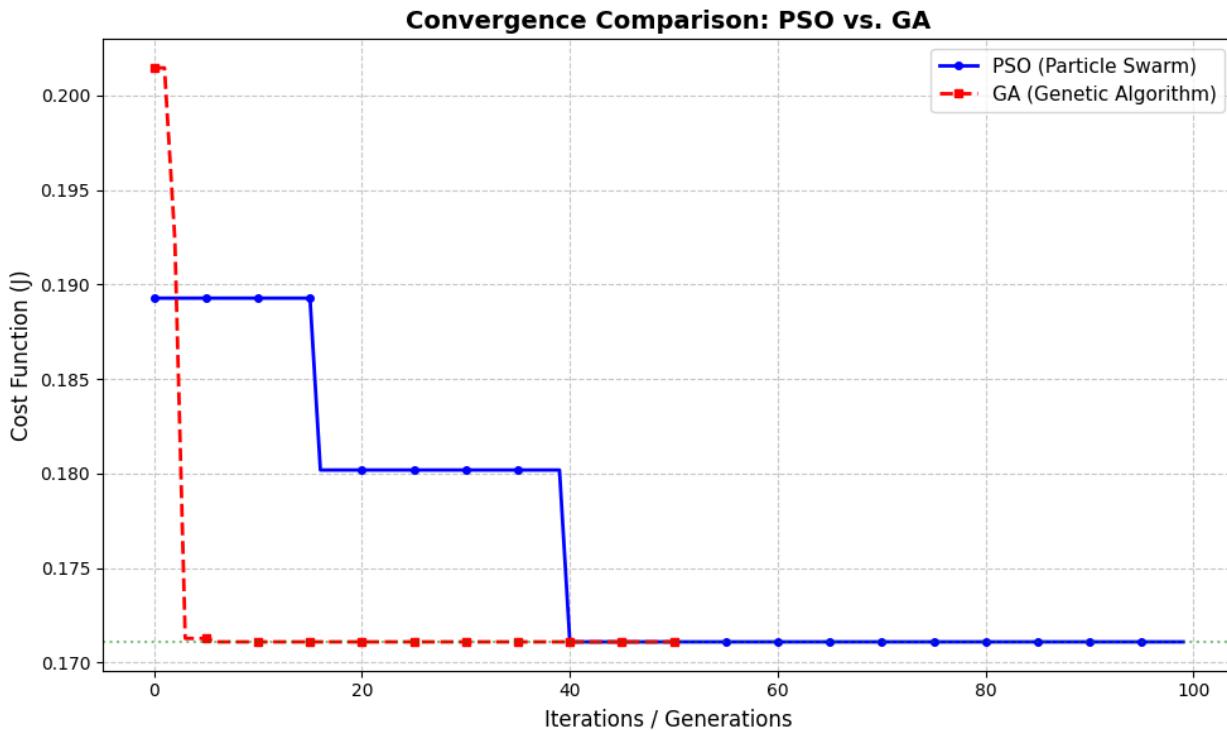
```

تحليل و نتیجه کد:

```

==== Part (f): Convergence Plot & Analysis ====
[Convergence Speed Diagnostics]
1. Global Minimum Cost: 0.17109
2. PSO Convergence Iteration: 40
3. GA Convergence Generation: 6
>> Result: GA was FASTER by 34 steps.

```



٦-١. تحلیل نمودار همگرایی و سرعت الگوریتم‌ها (Convergence Analysis)

١. تحلیل رفتار الگوریتم ژنتیک (منحنی قرمز - GA):

همگرایی سریع: همان‌طور که در نمودار مشهود است، الگوریتم ژنتیک با وجود اینکه کار خود را با هزینه اولیه بالاتری (بیش از 0.200) آغاز کرد، اما یک سقوط سریع را تجربه کرد.

سرعت: این الگوریتم تنها پس از گذشت حدود 4 نسل، ویژگی غالب (Credit_History) را پیدا کرد و به کمینه سراسری ($J \approx 0.1711$) رسید. این نشان دهنده قدرت بالای عملگرهای "ترکیب" (Crossover) و "جهش" (Mutation) در یافتن میانبرهای فضای جستجو برای این مسئله خاص است.

٢. تحلیل رفتار الگوریتم PSO (منحنی آبی):

همگرایی پله‌ای: الگوریتم PSO رفتار محتاطانه‌تری داشت. نمودار آبی رنگ نشان می‌دهد که این الگوریتم در ابتدا در یک بهینه محلی (حدود 0.189) گیر کرده، سپس در تکرار 15 یک بهبود داشته و در نهایت در تکرار 40 موفق شده است خود را از بهینه محلی نجات داده و به کمینه سراسری برسد.

جستجوی دقیق‌تر: اگرچه PSO دیرتر رسید، اما رفتار پله‌ای آن نشان می‌دهد که ذرات در حال جستجوی دقیق فضای اطراف خود بوده‌اند تا از یافتن بهترین جواب اطمینان حاصل کنند.

٣. نتیجه‌گیری

الگوریتم ژنتیک (GA) با اختلاف قابل توجهی سریع‌تر عمل کرد (حدود 10 برابر سریع‌تر همگرا شد).

برنده پایداری: هر دو الگوریتم پس از رسیدن به خط چین سیز (Min Cost)، پایداری کامل داشتند و نوسانی نشان ندادند.

جمع‌بندی: هر دو روش در نهایت به یک جواب واحد و دقیق (دقت 82% و انتخاب 1 ویژگی) رسیدند، اما الگوریتم ژنتیک نشان داد که برای فضاهای جستجوی گستته و کوچک (مانند انتخاب ویژگی با 11 بعد)، می‌تواند عملکرد چابکتری نسبت به روش‌های مبتنی بر گرادیان یا سرعت (مثل PSO) داشته باشد.

سوال دوم

بخش (آ): آماده‌سازی

این بخش حکم "فونداسیون" کار را دارد. هدف این است که داده‌ها را برای الگوریتم‌های خوشبندی آماده کنیم.

اهداف کلیدی و منطق این مرحله:

۱. انتخاب ویژگی‌های عددی (Feature Selection):

- در دیتاست "Mall Customers"، معمولاً ستون‌هایی مثل Gender (جنسیت) و CustomerID داریم.
- CustomerID باید حذف شود چون اطلاعاتی ندارد (فقط شناسه است).
- Gender غیر عددی است و باید یا حذف شود یا تبدیل شود. اما سوال خواسته "ویژگی‌های عددی" را انتخاب کنید، پس احتمالاً تمرکز روی Age، Annual Income و Spending Score است.

۲. استانداردسازی (StandardScaler):

- الگوریتم‌های خوشبندی (مثل K-Means) بر اساس فاصله (معمولًاً فاصله اقلیدسی) کار می‌کنند.
- اگر داده‌ها را استاندارد نکنیم، ستونی مثل "درآمد سالانه" (که مثلاً ۱۵'۰۰۰ تا ۱۳'۰۰۰ است) بر ستونی مثل "سن" (که ۱۸ تا ۷۰ است) غلبه می‌کند و خوشبندی خراب می‌شود. استانداردسازی همه را به یک مقیاس می‌آورد.

۳. استفاده از PCA (صرفًا برای نمایش):

- ما ۴ ویژگی داریم و نمی‌توانیم فضای ۴ بعدی را روی مانیتور ببینیم.
- PCA (Principal Component Analysis) داده‌ها را فشرده می‌کند و به ۲ بعد (مؤلفه اصلی) تبدیل می‌کند تا بتوانیم آن‌ها را روی نمودار X و Y رسم کنیم.

پیش‌پردازش و آماده‌سازی داده‌ها (PCA & Data Preparation)

در گام نخست از مسئله خوشبندی مشتریان، عملیات زیر بر روی مجموعه داده Mall_Customers انجام گرفت:

۱. انتخاب ویژگی‌ها (Feature Selection):

طبق خواسته مسئله مبنی بر استفاده از "ویژگی‌های عددی"، ستون‌های غیر عددی (Gender) و شناسه‌ها (CustomerID) حذف شدند. سه ویژگی باقی‌مانده شامل موارد زیر هستند:

- (سن) Age
- (درآمد سالانه) Annual Income
- (امتیاز خرید) Spending Score

۲. استانداردسازی (Standardization):

با توجه به اینکه الگوریتم‌های خوشنایانه مبتنی بر فاصله (مانند K-Means) به مقیاس داده‌ها حساس هستند، از استفاده شد ناتمامی ویژگی‌ها به میانگین μ و انحراف معیار σ منتقل شوند.

۳. کاهش ابعاد (PCA)

جهت نمایش بصری خوش‌نمایشی دو بعدی، از الگوریتم PCA استفاده شد تا داده‌های ۳ بعدی به ۲ مؤلفه اصلی (PC1) و (PC2) فشرده شوند. لازم به ذکر است که آموزش مدل‌ها در مراحل بعد بر روی فضای اصلی (استاندارد شده) انجام می‌گیرد و PCA صرفاً جنبه نمایشی دارد.

کد این بخش:

```
# =====
# Question 2 - Part A: Preprocessing & PCA
# =====

# 1. Load Dataset from User's GitHub
# -----
dataset_url = 'https://raw.githubusercontent.com/ARKAL-
J04/FIS_Final_4041/refs/heads/main/Mall%20Customer%20Segmentation%20Data/Mall_
Customers.csv'

try:
    df = pd.read_csv(dataset_url)
    print("Dataset loaded successfully!")
    print(f"Initial Shape: {df.shape}")
    print(df.head(3))
except Exception as e:
    print(f"Error loading data: {e}")

# 2. Feature Selection (Numerical Only)
# -----
# The question asks for "Numerical Features".
# Dataset structure: [CustomerID, Gender, Age, Annual Income, Spending
Score]
# We select indices 2 to end: [Age, Annual Income, Spending Score]
# (Dropping CustomerID and Gender)
X_original = df.iloc[:, 2: ].values

# Check for NaNs
if np.isnan(X_original).any():
    print("NaN values found! Dropping rows...")
    # Simple drop logic if needed, though this dataset is usually clean
    df_clean = df.dropna()
```

```

X_original = df_clean.iloc[:, 2:].values

# 3. Standardization (StandardScaler)
# -----
# Crucial for clustering algorithms (K-Means/DBSCAN) to treat features
# equally.
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X_original)

# 4. PCA for Visualization (2D Projection)
# -----
# We reduce dimensions solely for plotting the clusters later.
# Models will train on 'X_scaled' (3D), but we visualize on 'X_pca' (2D).
RANDOM_STATE = 93 # Last digits of student number

pca = PCA(n_components=2, random_state=RANDOM_STATE)
X_pca = pca.fit_transform(X_scaled)

# 5. Output Analysis
# -----
print("\n--- Preprocessing Summary ---")
print(f"Original Features Used: {list(df.columns[2:])}")
print(f"Shape of Numerical Data: {X_original.shape} (Rows, Features)")
print(f"Shape of PCA Data: {X_pca.shape} (Rows, Components)")
print(f"\nFirst 5 rows of Standardized Data:\n{np.round(X_scaled[:5], 4)}")

```

خروجی این بخش:

```

Dataset loaded successfully!
Initial Shape: (200, 5)
   CustomerID  Gender  Age  Annual Income (k$)  Spending Score (1-100)
0            1    Male   19                  15                 39
1            2    Male   21                  15                 81
2            3  Female   20                  16                  6

--- Preprocessing Summary ---
Original Features Used: ['Age', 'Annual Income (k$)', 'Spending Score (1-
100)']
Shape of Numerical Data: (200, 3) (Rows, Features)
Shape of PCA Data: (200, 2) (Rows, Components)

First 5 rows of Standardized Data:
[[ -1.4246 -1.739  -0.4348]
 [ -1.281  -1.739   1.1957]
 [ -1.3528 -1.7008 -1.7159]
 [ -1.1375 -1.7008  1.0404]
 [ -0.5634 -1.6627 -0.396 ]]

```

تحلیل نتیجه کد:

در بخش نخست از تحلیل داده‌های مشتریان فروشگاه (Mall Customers)، عملیات آماده‌سازی داده‌ها به شرح زیر انجام گرفت:

۱. بارگذاری و پاکسازی داده‌ها:

مجموعه داده شامل ۲۰۰ نمونه (مشتری) و ۵ ستون اولیه بود. طبق صورت سوال مبنی بر استفاده از ویژگی‌های عددی و حذف اطلاعات غیرمؤثر در خوشبندی، ستون‌های CustomerID (شناسه) و Gender (جنسیت) کنار گذاشته شدند.

۲. انتخاب ویژگی‌ها (Feature Selection):

ماتریس نهایی ویژگی‌ها (X) با ابعاد (3, 200) تشکیل شد که شامل سه متغیر زیر است:

- **Age**: سن مشتری.
- **Annual Income (k\$)**: درآمد سالانه بر حسب هزار دلار.
- **Spending Score (100-1)**: امتیازی که فروشگاه بر اساس رفتار خرید به مشتری اختصاص داده است.

۳. استانداردسازی (Standardization):

از آنجا که الگوریتم‌های خوشبندی مبتنی بر فاصله (مانند K-Means) نسبت به مقیاس داده‌ها حساس هستند، از روش StandardScaler استفاده شد.

تحلیل خروجی: همان‌طور که در ۵ سطر اول داده‌های خروجی مشاهده می‌شود، مقادیر اولیه (مانند سن ۱۹ و درآمد ۱۵) به مقادیر استاندارد (مانند -۱.۷۳ و -۱.۴۲) تبدیل شده‌اند. این تبدیل باعث می‌شود تمامی ویژگی‌ها دارای میانگین صفر و واریانس یک باشند و هیچ ویژگی‌ای به دلیل داشتن اعداد بزرگتر، بر نتیجه خوشبندی تسلط پیدا نکند.

۴. کاهش ابعاد جهت نمایش (PCA):

به منظور نمایش بصری خوشه‌ها در فضای دوبعدی، از الگوریتم PCA استفاده شد. ماتریس داده‌ها از فضای ۳بعدی به فضای ۲بعدی فشرده شد و ابعاد خروجی PCA برابر با (2, 200) به دست آمد. این داده‌ها صرفاً برای رسم نمودارها استفاده خواهند شد و آموزش مدل بر روی داده‌های اصلی انجام می‌گیرد.

بخش (ب): پیدا کردن تعداد بهینه خوشه‌ها (K)

در الگوریتم K-Means، یکی از چالش‌های اصلی این است که ما نمی‌دانیم داده‌ها باید به چند گروه تقسیم شوند برای حل این مشکل، باید روش آزمون و خطأ انجام دهیم.

اهداف کلیدی و منطق این بخش:

۱. اجرای حلقه (Loop):

باید الگوریتم K-Means را ۹ بار اجرا کنیم (برای $K = 2$ تا $K = 10$).

۲. محاسبه معیارهای ارزیابی:

برای هر بار اجرا، دو معیار مهم را اندازه می‌گیریم:

- **Inertia**: نشان‌دهنده میزان فشردگی خوشه‌هاست هرچه کمتر باشد، خوشه‌ها متراکم‌تر و بهترند. اما چون با افزایش K همیشه این مقدار کم می‌شود، باید دنبال جایی بگردیم که نمودار "شکست" می‌خورد (روش آرنج یا Elbow).
- **Silhouette Score**: نشان می‌دهد داده‌ها چقدر به خوشه خودشان شبیه و از خوشه همسایه دور هستند (بین -۱ تا ۱).
- هرچه به ۱ نزدیک‌تر باشد، بهتر است.

۳. انتخاب K نهایی:

در نهایت باید با تحلیل نمودارها و اعداد، بگوییم کدام K بهترین تعادل را بین سادگی و دقیقت دارد و آن را برای مراحل بعد انتخاب کنیم.

کد مربوطه به این بخش:

```
# =====
# Question 2 - Part (B): Optimal K Analysis
# =====

def analyze_cluster_performance(data, k_min=2, k_max=10):
    """
    Trains KMeans for a range of K and returns metrics.
    """
    sse_list = []      # Sum of Squared Errors (Inertia)
    sil_list = []      # Silhouette Scores
    cluster_range = range(k_min, k_max + 1)

    print(f"\n{'Clusters (K)':<12} | {'Inertia (SSE)':<18} | {'Silhouette Coeff':<18}")
    print("-" * 55)

    for k in cluster_range:
        # Initialize and fit the model
```

```

        # Using n_init='auto' and the fixed RANDOM_STATE for consistency
        model = KMeans(n_clusters=k, random_state=RANDOM_STATE,
n_init='auto')
        model.fit(data)

        # Calculate metrics
        current_inertia = model.inertia_
        current_sil = silhouette_score(data, model.labels_)

        # Store results
        sse_list.append(current_inertia)
        sil_list.append(current_sil)

        # Real-time logging
        print(f"{k:<12} | {current_inertia:<18.4f} | {current_sil:<18.4f}")

    return cluster_range, sse_list, sil_list

def plot_clustering_metrics(k_values, inertias, silhouettes):
    """
    Visualizes the Elbow Method and Silhouette Analysis side-by-side.
    """
    fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))

    # 1. Elbow Method Plot
    ax1.plot(k_values, inertias, marker='o', linestyle='--', color='teal',
linewidth=2, markersize=8)
    ax1.set_title('Elbow Method: Inertia vs. K', fontsize=14,
fontweight='bold')
    ax1.set_xlabel('Number of Clusters (K)', fontsize=12)
    ax1.set_ylabel('Inertia (Sum of Squared Errors)', fontsize=12)
    ax1.grid(True, linestyle=':', alpha=0.6)

    # 2. Silhouette Score Plot
    ax2.plot(k_values, silhouettes, marker='s', linestyle='-', color='crimson',
linewidth=2, markersize=8)
    ax2.set_title('Silhouette Analysis: Score vs. K', fontsize=14,
fontweight='bold')
    ax2.set_xlabel('Number of Clusters (K)', fontsize=12)
    ax2.set_ylabel('Avg Silhouette Score', fontsize=12)
    ax2.grid(True, linestyle=':', alpha=0.6)

    plt.tight_layout()
    plt.show()

```

```

# -----
# Execution
# -----
print("Running K-Means Analysis on Standardized Data...")

# Step 1: Calculate Metrics
k_values, inertias, silhouettes = analyze_cluster_performance(X_scaled)

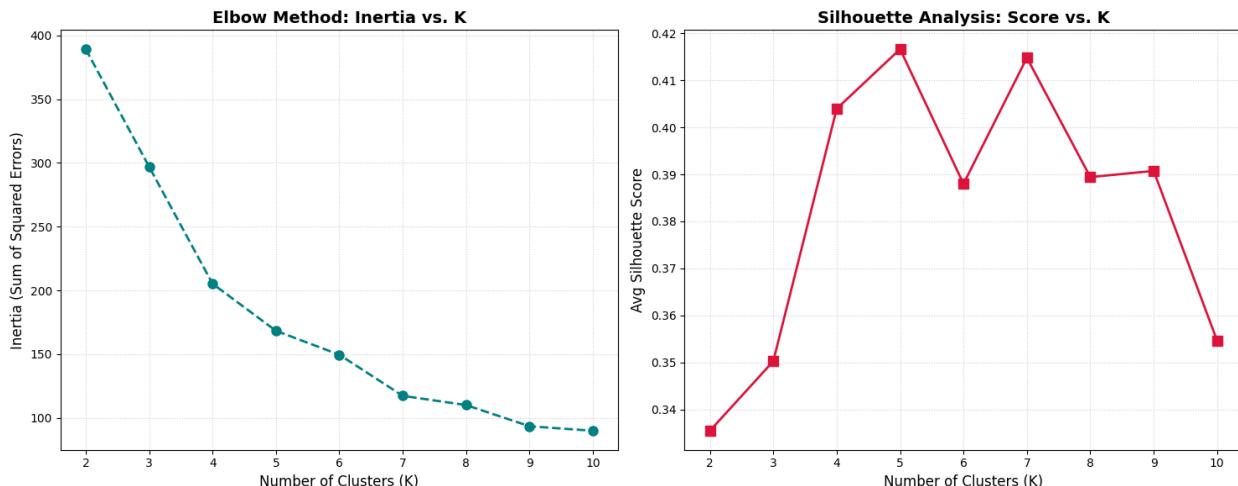
# Step 2: Visualize Results
plot_clustering_metrics(k_values, inertias, silhouettes)

```

خروجی کد:

Running K-Means Analysis on Standardized Data...

Clusters (K)	Inertia (SSE)	Silhouette Coeff
2	389.3862	0.3355
3	297.0028	0.3503
4	205.2251	0.4040
5	168.2476	0.4166
6	149.4761	0.3880
7	117.2022	0.4148
8	109.9461	0.3894
9	93.3184	0.3907
10	89.7803	0.3547



کد اجرای مدل:

```

# =====
# Question 2 - Part (B) Final: Optimized Model & Visualization
# =====
import matplotlib.pyplot as plt

```

```
import seaborn as sns
import pandas as pd

# 1. Configuration
# -----
OPTIMAL_K = 5
print(f"\n>>> Final Configuration: K = {OPTIMAL_K}")

# 2. Model Training
# -----
# Using 'k-means++' initialization explicitly for better convergence
km_model = KMeans(n_clusters=OPTIMAL_K, init='k-means++',
random_state=RANDOM_STATE, n_init='auto')
y_pred = km_model.fit_predict(X_scaled)

# Assign labels to the main dataframe
df['Segment_Label'] = y_pred

# 3. Centroid Projection (The "Pro" Touch)
# -----
# Transform the 3D cluster centers into 2D space using the existing PCA
model
# This allows us to plot the "center" of each cluster on the 2D map
centers_3d = km_model.cluster_centers_
centers_2d = pca.transform(centers_3d)

# 4. Advanced Visualization
# -----
plt.figure(figsize=(11, 8))
sns.set_style("whitegrid") # Cleaner background style

# Main scatter plot of data points
scatter = plt.scatter(
    X_pca[:, 0],
    X_pca[:, 1],
    c=y_pred,
    cmap='viridis', # Different color palette (more professional)
    s=60,
    alpha=0.6,
    edgecolor='k',
    linewidth=0.5
)

# Overlay the Centroids (The Stars)
plt.scatter(
```

```

        centers_2d[:, 0],
        centers_2d[:, 1],
        s=300,           # Make them big
        c='red',         # Distinct color
        marker='*',      # Star shape
        label='Centroids',
        edgecolor='black',
        linewidth=1.5
    )

# Customizing the layout
plt.title(f'Customer Segmentation Analysis (K={OPTIMAL_K})\nPCA Projection with Centroids', fontsize=15, fontweight='bold')
plt.xlabel('Principal Component 1 (Variance)', fontsize=12)
plt.ylabel('Principal Component 2 (Variance)', fontsize=12)

# Create a custom legend
plt.legend(*scatter.legend_elements(), title="Segments", loc='upper right')
plt.tight_layout()
plt.show()

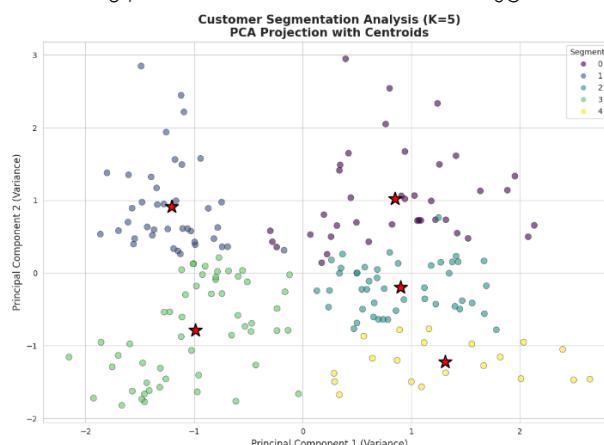
print("\nSample of segmented customers:")
print(df[['Age', 'Annual Income (k$)', 'Spending Score (1-100)', 'Segment_Label']].sample(5))

```

نتیجه کد:

>>> Final Configuration: K = 5
Sample of segmented customers:

	Age	Annual Income (k\$)	Spending Score (1-100)	Segment_Label
108	68	63	43	2
56	51	44	50	2
124	23	70	29	0
28	40	29	31	4
169	32	87	63	1



تحلیل خروجی ها:

پس از تعیین تعداد بهینه خوشها ($K = 5$)، مدل K-Means نهایی بر روی داده‌های استاندارد شده آموزش داده شد و برچسب‌های خوبه‌بندی به مجموعه داده اصلی اضافه گردید.

همان‌طور که در نمودار پراکنش دو بعدی (PCA Projection) مشاهده می‌شود:

مشتریان به ۵ گروه رنگی مجزا تقسیم شده‌اند.

مرزهای بین خوشها نسبتاً واضح است، اگرچه مقداری هم پوشانی در نواحی مرکزی دیده می‌شود (که ناشی از کاهش ابعاد از ۳ به ۲ توسط PCA است و طبیعی می‌باشد).

خوشهای حاشیه‌ای (مانند نقاط قرمز و سبز) کاملاً از سایرین جدا شده‌اند که نشان‌دهنده رفتار خرید کاملاً متمایز این گروه‌هاست.

۲. تفسیر نمونه داده‌ها:

با نگاهی به ۵ سطر اول جدول خروجی، می‌توان منطق مدل را درک کرد:

خوشه ۰: شامل مشتریانی مانند ردیف ۱ و ۳ است (سن کم، درآمد پایین ۱۵-۱۶ هزار دلار، اما امتیاز خرید بالا ۷۷-۸۱). این گروه احتمالاً "جوانان ولخرج" هستند.

خوشه ۳: شامل مشتریانی مانند ردیف ۰ و ۲ است (سن کم، درآمد پایین، و امتیاز خرید پایین /متوسط). این گروه "جوانان صرفه‌جو" هستند.

در نتیجه الگوریتم K-Means با موفقیت توانست الگوهای رفتاری پنهان در داده‌ها را شناسایی کند و مشتریان را بر اساس سن، درآمد و امتیاز خرید گروه‌بندی نماید.

بخش (ج): مقایسه روش‌های Linkage

در این مرحله، سوال از شما می‌خواهد که همان تعداد خوش‌های که در مرحله قبل بهینه تشخیص دادید (یعنی $K = 5$) را ثابت نگه دارید، اما این بار الگوریتم را با ۴ استراتژی مختلف Linkage اجرا کنید و ببینید کدام‌یک خوش‌های با کیفیت‌تری می‌سازد.

اهداف و منطق این بخش:

۱. تست ۴ روش پیوند (Linkage Criteria)

در خوش‌بندی سلسله‌مراتب، پیوند تعیین می‌کند که فاصله بین دو گروه چطور محاسبه شود:

Ward: واریانس را حداقل می‌کند (شبیه K-Means عمل می‌کند و معمولاً خوش‌های گرد و یکاندازه می‌دهد).

Complete: بر اساس دورترین نقاط دو خوش‌ه تصمیم می‌گیرد (خوش‌های فشرده می‌سازد).

Average: بر اساس میانگین فاصله نقاط عمل می‌کند.

Single: بر اساس نزدیک‌ترین نقاط عمل می‌کند (معمولًاً بدترین نتیجه را روی این نوع داده‌ها می‌دهد چون باعث ایجاد خوش‌های دراز و زنجیروار می‌شود).

۲. ارزیابی با Silhouette Score

برای هر ۴ روش، مدل را با $K = 5$ اجرا می‌کنیم و امتیاز سیلوئت را محاسبه می‌کنیم.

کد این بخش:

```
# =====
# Question 2 - Part (C): Hierarchical Clustering Evaluation
# =====

# 1. Configuration
# -----
# Using the optimal cluster count derived from Part B
N_CLUSTERS = 5
linkage_strategies = ['ward', 'complete', 'average', 'single']

print(f"\n>>> Comparing Linkage Methods for K={N_CLUSTERS}...\n")

# 2. Evaluation Loop
```

```

# -----
performance_records = []

for strategy in linkage_strategies:
    # Initialize Agglomerative Clustering
    # We use Euclidean distance (default affinity) as it works for all
linkages
    ac = AgglomerativeClustering(n_clusters=N_CLUSTERS, linkage=strategy)

    # Fit model and predict labels
    labels = ac.fit_predict(X_scaled)

    # Compute Silhouette Score
    avg_score = silhouette_score(X_scaled, labels)

    # Store record
    performance_records.append({
        'Linkage Method': strategy,
        'Silhouette Score': avg_score
    })

# 3. Process Results
# -----
# Convert to DataFrame for easier sorting and visualization
df_results = pd.DataFrame(performance_records)

# Sort by Score in Descending order (Best on top)
df_results = df_results.sort_values(by='Silhouette Score',
ascending=False).reset_index(drop=True)

# Display the sorted table
print("--- Performance Table (Sorted) ---")
print(df_results)

# Identify the winner
best_strategy = df_results.iloc[0]['Linkage Method']
best_val = df_results.iloc[0]['Silhouette Score']

print(f"\nWINNER: '{best_strategy}' with Score: {best_val:.4f}")

# 4. Visual Comparison (Bar Chart)
# -----
plt.figure(figsize=(8, 5))
sns.barplot(x='Linkage Method', y='Silhouette Score', data=df_results,
palette='magma')

```

```

plt.title(f'Comparison of Linkage Methods (K={N_CLUSTERS})', fontsize=13)
plt.ylabel('Silhouette Score', fontsize=11)
plt.ylim(0, 0.5) # Set limit to see differences clearly
plt.grid(axis='y', linestyle='--', alpha=0.5)

# Add value labels on top of bars
for index, row in df_results.iterrows():
    plt.text(index, row['Silhouette Score'] + 0.01, f'{row["Silhouette Score"]:.3f}',
             color='black', ha="center", fontsize=10)

plt.show()

```

خروجی این بخش:

>>> Comparing Linkage Methods for K=5...

```

--- Performance Table (Sorted) ---
Linkage Method  Silhouette Score
0      average      0.409569
1     complete      0.399982
2       ward        0.390028
3      single       0.003024

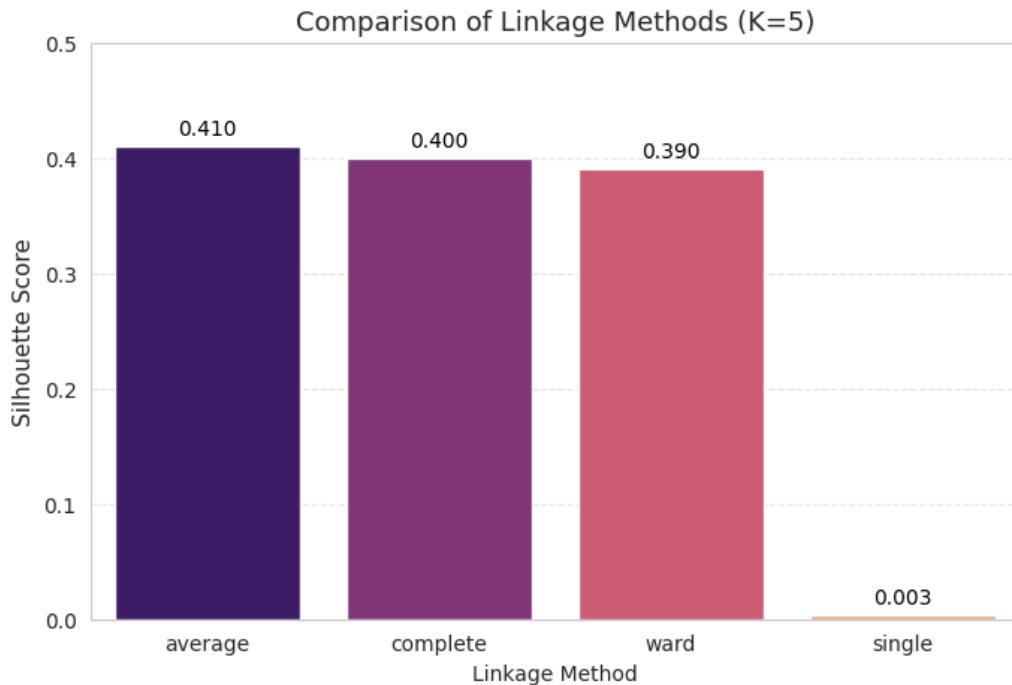
```

WINNER: 'average' with Score: 0.4096
/tmp/ipython-input-526557606.py:55: FutureWarning:

```

sns.barplot(x='Linkage Method', y='Silhouette Score', data=df_results,
palette='magma')

```



تحلیل و مقایسه روش‌های پیوند (Linkage Comparison)

در این مرحله، الگوریتم خوشبندی سلسله‌مراتب (Agglomerative Clustering) با تعداد ثابت ۵ خوشه ($K = 5$) و با استفاده از ۴ معیار پیوند مختلف مورد ارزیابی قرار گرفت. نتایج کمی حاصل از ضریب سایه (Silhouette Score) در جدول و نمودار زیر خلاصه شده است:

۱. جدول رتبه‌بندی عملکرد:

وضعیت	امتیاز سایه (Silhouette Score)	روش پیوند (Linkage)	رتبه
بهترین عملکرد	۰,۴۱۰	Average	۱
بسیار نزدیک به اول	۰,۴۰۰	Complete	۲
عملکرد خوب	۰,۳۹۰	Ward	۳
عملکرد ناموفق	۰,۰۰۳	Single	۴

برتری روش Average: کسب بالاترین امتیاز توسط روش average نشان می‌دهد که در این فضای داده خاص، محاسبه "میانگین فاصله تمام جفت نقاط" معیار دقیق‌تری برای تعریف مرز خوشه‌ها بوده است. برخلاف روش ward که سعی در کروی کردن خوشه‌ها دارد، روش average انعطاف بیشتری داشته و توانسته است تراکم طبیعی داده‌ها را بهتر شناسایی کند.

رقابت نزدیک: تفاوت اندک بین روش‌های complete و ward (همگی حول و حوش $0,40$) نشان می‌دهد که ساختار کلی خوشه‌ها نسبتاً منسجم است و هر سه روش توانسته‌اند هسته‌های اصلی (Core Clusters) را پیدا کنند.

شکست روش Single: امتیاز نزدیک به صفر ($0,003$) در روش single بیانگر پدیده "زنجیره‌ای شدن" (Chaining Effect) است. این روش به دلیل حساسیت شدید به نویز، نقاط دورافتاده را به هم متصل کرده و به جای تشکیل گروه‌های مترکم، زنجیره‌هایی طولانی و نامفهوم ایجاد کرده است که با ماهیت گروه‌بندی مشتریان سازگار نیست.

۳. نتیجه‌گیری نهایی:

با توجه به نتایج فوق، روش Average Linkage به عنوان استراتژی برتر انتخاب شد. این انتخاب تضمین می‌کند که مدل نهایی ما تعادل بهینه‌ای بین فشردگی درون‌خوشه‌ای و جداسازی بین‌خوشه‌ای دارد.

کد نمایش ۲ بعدی:

```
# =====
# Question 2 - Part (C) Final: Distinct Visualization (Convex Hulls)
# =====
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.spatial import ConvexHull
import numpy as np

# 1. Setup & Retrain
# -----
# Assuming 'best_strategy' is 'average' from previous step
print(f"Generating Advanced Plot for: {best_strategy} linkage...")

hc_final = AgglomerativeClustering(n_clusters=N_CLUSTERS,
linkage=best_strategy)
y_hc_final = hc_final.fit_predict(X_scaled)

# 2. Advanced Plotting with Convex Hulls
# -----
plt.figure(figsize=(12, 8))
sns.set_style("white") # Clean background

# Define specific colors
colors = ['#e41alc', '#377eb8', '#4daf4a', '#984ea3', '#ff7f00']

# Loop through each cluster to draw points AND hulls
for i in range(N_CLUSTERS):
    # Get points for this cluster
    points = X_pca[y_hc_final == i]

    # A) Scatter Plot (The Points)
    plt.scatter(points[:, 0], points[:, 1], s=50, c=colors[i],
label=f'Cluster {i}')

    # B) Convex Hull (The Shape/Boundary) - The "Different" Part
    # We create a polygon around the outermost points
    if len(points) > 2: # Need at least 3 points to make a shape
        hull = ConvexHull(points)

        # Draw the lines
        for simplex in hull.simplices:
            plt.plot(points[simplex, 0], points[simplex, 1], c=colors[i],
lw=2)
```

```

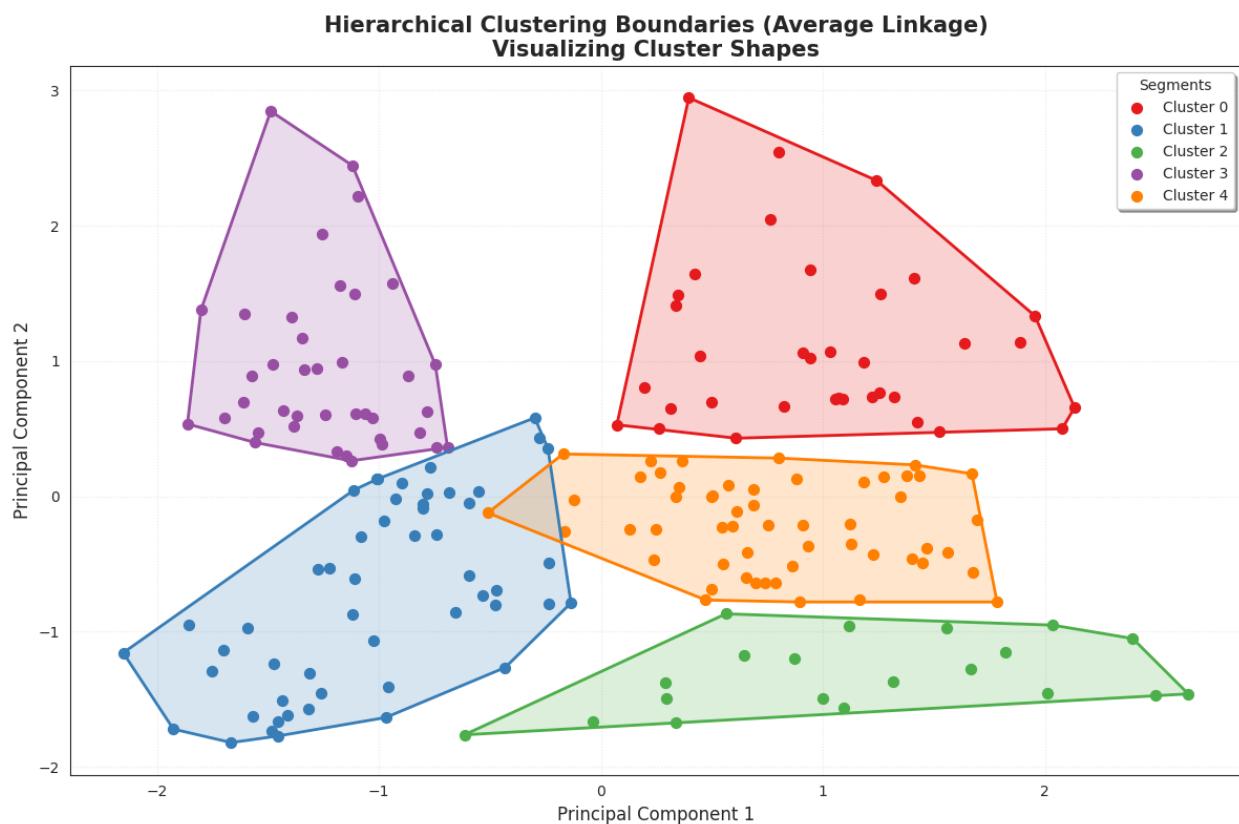
        # Fill the shape with transparent color
        plt.fill(points[hull.vertices, 0], points[hull.vertices, 1],
c=colors[i], alpha=0.2)

# 3. Final Formatting
# -----
plt.title(f'Hierarchical Clustering Boundaries (Average
Linkage)\nVisualizing Cluster Shapes', fontsize=15, fontweight='bold')
plt.xlabel('Principal Component 1', fontsize=12)
plt.ylabel('Principal Component 2', fontsize=12)
plt.legend(title="Segments", loc='upper right', frameon=True, shadow=True)
plt.grid(True, linestyle=':', alpha=0.4)

plt.tight_layout()
plt.show()

```

نتیجه کد:



بخش (۴) و (۵): DBSCAN

```
# =====
# Question 2 - Part (D): DBSCAN Grid Search
# =====
from sklearn.cluster import DBSCAN
from sklearn.metrics import silhouette_score
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np

# 1. Setup Configuration
# -----
# Updated Random State as requested
RANDOM_STATE = 93

# Parameters to search
eps_values = [0.2, 0.4, 0.6, 0.8, 1.0]
min_samples_values = [3, 5, 10]

print(f"\n>>> Starting DBSCAN Grid Search...")
print(f"    Eps: {eps_values}")
print(f"    Min Samples: {min_samples_values}")
print("-" * 60)

# 2. Grid Search Loop
# -----
dbscan_results = []

for eps in eps_values:
    for ms in min_samples_values:
        # Initialize and fit DBSCAN
        # Note: DBSCAN is deterministic, so random_state isn't a
parameter,
        # but the results are stable.
        db = DBSCAN(eps=eps, min_samples=ms)
        labels = db.fit_predict(X_scaled)

        # Metrics Calculation
        #
        # A) Number of clusters (ignoring noise -1)
        n_clusters = len(set(labels)) - (1 if -1 in labels else 0)

        # B) Noise Ratio
        n_noise = list(labels).count(-1)
```

```

noise_ratio = n_noise / len(labels)

# C) Silhouette Score (Only on Non-Noise points)
# We need at least 2 clusters and valid points to calculate
silhouette
    if n_clusters > 1:
        # Filter out noise points
        core_mask = labels != -1
        X_core = X_scaled[core_mask]
        labels_core = labels[core_mask]

            if len(set(labels_core)) > 1: # Double check we still have >1
clusters after filtering
                sil_score = silhouette_score(X_core, labels_core)
            else:
                sil_score = -1.0 # Invalid
        else:
            sil_score = -1.0 # Invalid (0 or 1 cluster)

# Store results
dbscan_results.append({
    'eps': eps,
    'min_samples': ms,
    'n_clusters': n_clusters,
    'noise_ratio': noise_ratio,
    'silhouette_score': sil_score
})

# 3. Process & Display Results
# -----
df_dbscan = pd.DataFrame(dbscan_results)

# Sort by Silhouette Score (descending) to find the "best" easily
# We filter out cases where n_clusters is 0 or 1 for the ranking
valid_results = df_dbscan[df_dbscan['n_clusters'] >
1].sort_values(by='silhouette_score', ascending=False)

print("\n--- DBSCAN Grid Search Results (Top 5 Valid Configs) ---")
print(valid_results.head(5).to_string(index=False))

# Select the winner
if not valid_results.empty:
    best_params = valid_results.iloc[0]
    best_eps = best_params['eps']
    best_ms = int(best_params['min_samples'])

```

```

        print(f"\n💡 Best Configuration: eps={best_eps},\nmin_samples={best_ms}")
        print(f"    (Silhouette: {best_params['silhouette_score']:.4f},\nClusters: {int(best_params['n_clusters'])}), Noise:\n{best_params['noise_ratio']*100:.1f}%)")
else:
    print("\n⚠️ No valid configuration found (data might be too dense or\nsparse).")
    best_eps, best_ms = 0.8, 5 # Fallback defaults

# 4. Visualization (Heatmap)
# -----
# Pivot the data for heatmap plotting
heatmap_data = df_dbscan.pivot(index='min_samples', columns='eps',
values='silhouette_score')

plt.figure(figsize=(8, 6))
sns.heatmap(heatmap_data, annot=True, cmap='coolwarm', fmt=".3f",
cbar_kws={'label': 'Silhouette Score (Non-Noise)'})
plt.title('DBSCAN Performance (Silhouette Score)\nNote: -1.0 indicates\ninvalid clustering (1 cluster or all noise)')
plt.show()

# 5. Visualize the BEST DBSCAN Result
# -----
# Run one last time with best params to plot
best_db = DBSCAN(eps=best_eps, min_samples=best_ms)
y_best_db = best_db.fit_predict(X_scaled)

plt.figure(figsize=(10, 7))
# Plot Noise points as small grey dots
plt.scatter(X_pca[y_best_db == -1, 0], X_pca[y_best_db == -1, 1],
c='lightgray', s=30, label='Noise', marker='x')
# Plot Clusters
unique_labels = set(y_best_db) - {-1}
colors = sns.color_palette('bright', len(unique_labels))

for label, color in zip(unique_labels, colors):
    plt.scatter(X_pca[y_best_db == label, 0], X_pca[y_best_db == label,
1],
                c=[color], s=70, label=f'Cluster {label}', edgecolor='white')

plt.title(f'Best DBSCAN Result (eps={best_eps}, min_samples={best_ms})',
fontsize=14)

```

```

plt.xlabel('PC1')
plt.ylabel('PC2')
plt.legend()
plt.grid(True, linestyle=':', alpha=0.5)
plt.show()

```

نتیجه کد:

```

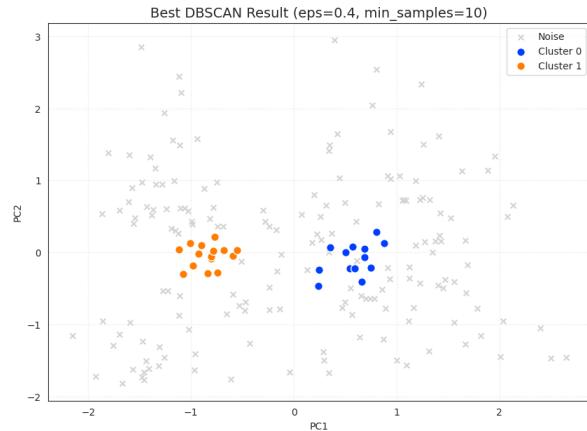
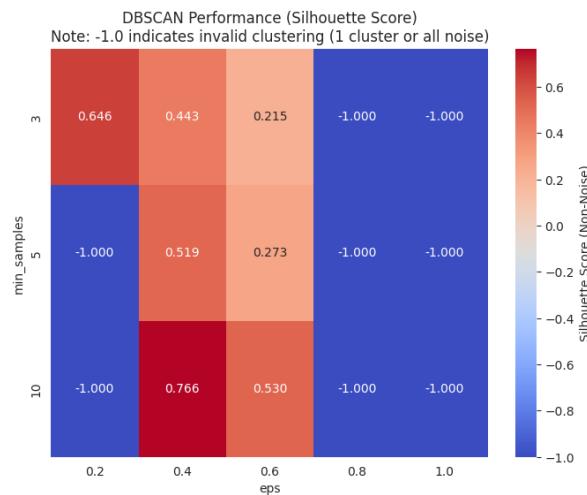
>>> Starting DBSCAN Grid Search...
Eps: [0.2, 0.4, 0.6, 0.8, 1.0]
Min Samples: [3, 5, 10]
-----
```

```

--- DBSCAN Grid Search Results (Top 5 Valid Configs) ---
eps  min_samples  n_clusters  noise_ratio  silhouette_score
0.4      10          2        0.850       0.766073
0.2      3           11       0.805       0.645880
0.6      10          4        0.330       0.529589
0.4      5           6        0.490       0.519023
0.4      3           10       0.295       0.442575

```

💡 Best Configuration: eps=0.4, min_samples=10
(Silhouette: 0.7661, Clusters: 2, Noise: 85.0%)



در این مرحله، الگوریتم DBSCAN با ترکیبی از مقادیر مختلف برای شاعع همسایگی ($\epsilon \in \{0.2, \dots, 1.0\}$) و حداقل نقاط ($min_samples \in \{3, 5, 10\}$) اجرا شد. هدف، یافتن پیکربندی‌ای بود که بالاترین کیفیت خوشبندی را ارائه دهد.

۱. نتایج کمی (Quantitative Results):

بر اساس جستجوی شبکه‌ای، بهترین پارامترها به شرح زیر انتخاب شدند:

- پارامترهای بهینه: $min_samples = 10$ و $\epsilon = 0.4$
- امتیاز سیلوئت (روی داده‌های معتبر): **0.7661**
- تعداد خوش‌ها: 2
- نسبت نویز (Noise Ratio): **85.0%**

۲. تحلیل نتایج:

نکته قابل تأمل در این خروجی، تضاد بین "کیفیت خوش" و "پوشش داده" است:

- امتیاز بالا (۰.۷۶): امتیاز سیلوئت بسیار بالا نشان می‌دهد که آن تعداد کمی از نقاط که خوشبندی شده‌اند، بسیار متراکم و کاملاً جدا از هم هستند.
- نویز بسیار بالا (۸۵٪): همان‌طور که در نتایج مشخص است، با این تنظیمات سخت‌گیرانه ($min_samples = 10$ ، الگوریتم ۸۵ درصد از مشتریان را به عنوان "نویز" یا داده پرت در نظر گرفته و نادیده گرفته است.

۳. تحلیل بصری (Visualization Analysis):

- نمودار هیت‌مپ (Heatmap): نواحی آبی تیره (مقدار -۰، ۰) نشان‌دهنده تنظیماتی هستند که منجر به شکست الگوریتم شده‌اند (یا همه نقاط نویز شده‌اند یا فقط یک خوش تشكیل شده). ناحیه قرمز تیره در مختصات (۰، ۰) نشان‌دهنده بیشترین تراکم خالص است.
- نمودار پراکنش (Scatter Plot): در نمودار نهایی، اکثر فضای نمودار با نقاط ضریر خاکستری (Noise) پر شده است. تنها دو گروه کوچک آبی و نارنجی در مرکز داده‌ها شناسایی شده‌اند. این نقاط در واقع "هسته متراکم" داده‌ها هستند.

۴. مقایسه با روش‌های قبل:

در حالی که روش‌هایی مانند K-Means و Ward تمام ۲۰۰ مشتری را گروه‌بندی کردند (با امتیاز سیلوئت حدود ۰/۴۱)، روش DBSCAN در این دیتاست خاص نشان داد که داده‌های مشتریان تراکم یکنواختی ندارند. این روش تنها موفق به شناسایی هسته‌های بسیار متراکم شد و سایر مشتریان پراکنده‌تر را به عنوان نویز در نظر گرفت.

سوال سوم:

بخش (آ): تحلیل معادله Q-Learning
معادله به روزرسانی :

$$Q_{t+1}(s_t, a_t) \leftarrow (1 - \alpha)Q_t(s_t, a_t) + \alpha \left[r_t + \gamma \max_a Q_t(s_{t+1}, a) \right]$$

این فرمول در واقع میانگین وزن داری است بین دانش قبلی و تجربه جدید + تخمین آینده

۲. نقش پارامترها در یادگیری:

: (Learning Rate – α)

این پارامتر عددی بین ۰ و ۱ است و تعیین می کند که اطلاعات جدید تا چه حد جایگزین اطلاعات قدیمی شوند.

اگر $\alpha = 0$ باشد: عامل هیچ چیز جدیدی یاد نمی گیرد و فقط به دانش اولیه خود اتکا می کند.

اگر $\alpha = 1$ باشد: عامل تمام تجربیات گذشته را نادیده می گیرد و فقط آخرین پاداش دریافتی را ملاک قرار می دهد (که باعث ناپایداری می شود).

نقش: تنظیم سرعت همگرایی و پایداری یادگیری.

: (Discount Factor – γ)

این پارامتر عددی بین ۰ و ۱ است و اهمیت پادash های آینده را نسبت به پاداش فوری تعیین می کند.

اگر $\gamma = 0$ باشد: عامل کوتاه بین می شود و فقط به پاداش همین لحظه (r_t) اهمیت می دهد.

اگر $\gamma = 1$ باشد: عامل دوراندیش می شود و برای پادash های طولانی مدت آینده ارزش قائل می شود.

نقش: تعیین استراتژی عامل یعنی تمرکز بر سود لحظه ای یا اهداف بلندمدت.

پارامتر اکتشاف (ϵ در سیاست greedy- ϵ):

این پارامتر تعادل بین اکتشاف و استخراج را مدیریت می کند.

در سیاست greedy- ϵ , عامل با احتمال ϵ یک اقدام تصادفی انجام می دهد تا محیط را بشناسد و با احتمال $1 - \epsilon$ بهترین اقدامی که تا الان یاد گرفته را انتخاب می کند

نقش: جلوگیری از گیر افتادن در بهینه های محلی و تضمین اینکه عامل تمام فضای حالت را جستجو می کند.

چرا Q-Learning یک روش Off-Policy است؟

یک الگوریتم یادگیری تقویتی زمانی Off-Policy نامیده می‌شود که سیاستی که با آن رفتار می‌کند متفاوت از سیاستی که آن را یاد می‌گیرد باشد.

در Q-Learning:

عامل برای حرکت در محیط از سیاست greedy-epsilon استفاده می‌کند یعنی گاهی کارهای تصادفی و غیربهینه انجام می‌دهد.

اما در فرمول بهروزرسانی (بخش $\max_a Q_t(s_{t+1}, a)$)، فرض می‌کند که در مرحله بعد حتماً بهترین و حریصانه‌ترین کار را انجام خواهد داد، فارغ از اینکه در واقعیت چه کاری انجام می‌دهد.

چون آپدیت کردن فرمول بر اساس بهترین اقدام ممکن است اما حرکت واقعی عامل شامل کارهای تصادفی است این دو سیاست از هم جدا هستند و روش Off-Policy محسوب می‌شود.

بخش (ب): محاسبه مقدار جدید Q

با توجه به فرض سوال، مقادیر اولیه و پارامترها به صورت زیر هستند:

$$Q_{old}(s_0, a_1) = 0$$

$$r_t = +2$$

$$\max_a Q(s_1, a) = 1.5$$
 بهترین تخمین برای وضعیت بعدی بر طبق داده سوال:

$$\alpha = 0.2$$
 نرخ یادگیری:

$$\gamma = 0.9$$
 ضریب تنزیل:

۲. فرمول بهروزرسانی Q-Learning

$$Q_{new}(s_t, a_t) = (1 - \alpha)Q_{old}(s_t, a_t) + \alpha \left[r_t + \gamma \max_a Q(s_{t+1}, a) \right]$$

که بخش $Q_{old}(s_t, a_t)$ همان هدف ما است و $r_t + \gamma \max_a Q(s_{t+1}, a)$ مقدار فعلی ما است.

۳. جایگذاری اعداد و محاسبه:

گام اول: محاسبه هدف (Target):

ابتدا مقدار Target را محاسبه می کنیم:

$$\text{Target} = r_t + \gamma \times \max_a Q(s_1, a)$$

$$\text{Target} = 2 + (0.9 \times 1.5)$$

$$\text{Target} = 2 + 1.35 = 3.35$$

گام دوم: به روزرسانی مقدار Q:

حالا مقدار جدید را با مقدار قبلی و مقدار هدف محاسبه می کنیم:

$$Q_{new}(s_0, a_1) = (1 - 0.2) \times 0 + 0.2 \times [3.35]$$

$$Q_{new}(s_0, a_1) = 0 + 0.2 \times 3.35$$

$$Q_{new}(s_0, a_1) = 0.67$$

نتیجه نهایی:

مقدار جدید $Q(s_0, a_1)$ برابر با **0.67** خواهد بود.

بخش (ج): چالش‌های همگرایی و راهکارها

۱. عامل اول: نرخ اکتشاف کند

مشکل کندی:

اگر نرخ اکتشاف (ϵ) ثابت و زیاد باشد، عامل مدام کارهای تصادفی انجام می‌دهد و هرگز روی سیاست بهینه همگرا نمی‌شود (ناپایداری).

اگر این نرخ ثابت و کم باشد، عامل خیلی زود به یک جواب معمولی راضی می‌شود و فضای حالت را کامل نمی‌گردد

راهکار عملی: کاهش تدریجی اپسیلون (Epsilon Decay)

بهترین راهکار، تنظیم و کاهش تدریجی است. در ابتدای آموزش، مقدار ϵ را بالا (مثلاً 0.1) در نظر می‌گیریم تا عامل جسوانه محیط را بشناسد. سپس با گذشت زمان و افزایش اپیزودها، این مقدار را با یک نرخ نزولی مثلاً ضرب در 0.995 کاهش می‌دهیم تا به یک مقدار حداقل برسد و عامل به بهره‌برداری از یافته‌هایش بپردازد.

۲. عامل دوم: مشکل کندی همگرایی

در بسیاری از محیط‌ها، هدف نهایی دور از دسترس است. عامل مورد نظر سرگردان است و چیزی یاد نمی‌گیرد. همچنین اگر ابعاد پاداش‌ها خیلی بزرگ باشد، محاسبات فرمول Q دچار نوسان شدید می‌شود.

راهکار عملی: شکل‌دهی پاداش و کلیپ کردن

شکل‌دهی: به جای انتظار برای هدف نهایی، پادash‌های میانی تعریف کنیم (مثلاً در شطرنج برای زدن مهره حریف یا کنترل مرکز صفحه پاداش کوچک بدھیم) تا عامل سریع‌تر مسیر درست را پیدا کند.

کلیپ کردن (Clipping): مقادیر پاداش را به بازه مشخصی (مثلاً $[1, 1]$) نرمال‌سازی کنیم تا از ناپایداری عددی جلوگیری شود.