QUESTION # 01:

Importing necessary modules:

```
import pandas as pd
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
import numpy as np
from sklearn.metrics import classification_report
```

EVALUATION PARAMETERS:

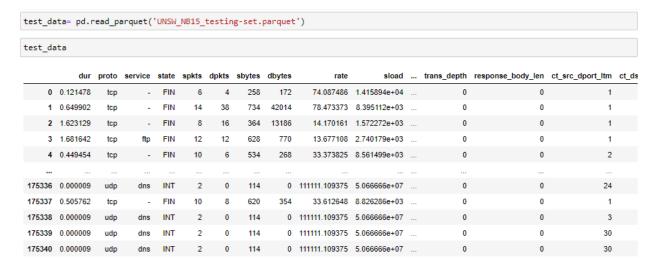
1. DATASET:

UNSW-NB15 dataset contains network traffic data with labeled intrusion instances.

Reading training and testing data

rain_	data= pd	l.read_	_parquet	t('UNS	W_NB1	5_trai	ning-se	et.parq	uet')					
rain_	data													
	dur	proto	service	state	spkts	dpkts	sbytes	dbytes	rate	sload	 trans_depth	response_body_len	ct_src_dport_ltm	ct_d
0	0.000011	udp	-	INT	2	0	496	0	90909.093750	1.803636e+08	 0	0	1	
1	8000008	udp	-	INT	2	0	1762	0	125000.000000	8.810000e+08	 0	0	1	
2	0.000005	udp		INT	2	0	1068	0	200000.000000	8.544000e+08	 0	0	1	
3	0.000006	udp	-	INT	2	0	900	0	166666.656250	6.000000e+08	 0	0	2	
4	0.000010	udp	-	INT	2	0	2126	0	100000.000000	8.504000e+08	 0	0	2	
82327	0.000005	udp	-	INT	2	0	104	0	200000.000000	8.320000e+07	 0	0	1	
82328	1.106101	tcp	-	FIN	20	8	18062	354	24.410067	1.241044e+05	 0	0	1	
82329	0.000000	arp	-	INT	1	0	46	0	0.000000	0.000000e+00	 0	0	1	
82330	0.000000	arp	-	INT	1	0	46	0	0.000000	0.000000e+00	 0	0	1	
82331	0.000009	udp	-	INT	2	0	104	0	111111.109375	4.622222e+07	 0	0	1	

82332 rows × 36 columns



175341 rows × 36 columns

2. Preprocessing: Standardization of numerical features and one-hot encoding of categorical features.

Preprocessing:

```
extracting numerical columns
```

numerical_transformer = StandardScaler()

```
numeric_cols = train_data.select_dtypes(include=np.number).columns.tolist()

removing 'label' column::

if 'label' in numeric_cols:
    numeric_cols.remove('label')

print("Numerical Columns:", numeric_cols)

Numerical Columns: ['dur', 'sbytes', 'dbytes', 'spkts', 'dpkts', 'rate', 'sload', 'dload']

extracting categorical columns

categorical_cols = train_data.select_dtypes(exclude=np.number).columns.tolist()

print("Categorical Columns: ['proto']

Initializing a StandardScaler for numerical features
```

2. GENETIC ALGORITHM IMPLEMENTATION & PARAMETERS:

Applying genetic algorithm:

```
population_size = 100

chromosome_length = x_train_preprocessed.shape[1]
print(chromosome_length)

139

population = np.random.randint(0, 2, (population_size, chromosome_length))
print(population)

[[0 1 ... 1 0 1]
[[0 1 0 ... 0 0 0]
[[1 1 0 ... 1 0]
...
[[0 1 0 ... 0 0 0]
[[0 1 0 ... 1 0 0]
[[0 0 0 ... 1 0 0]]
```

- 1. **Population Size:** Experimented with population sizes of 50, 100, and 200.
- 2. Chromosome Length: Determined by the number of features after preprocessing.

3. **Fitness Function**: Custom fitness function based on true positives, true negatives, false positives, and false negatives.

function to calculate fitness

```
idef calculate_fitness(individual, features, labels):
    if isinstance(features, np.ndarray):
        features_dense = features
    else:
        features_dense = features.toarray()
    individual = np.array(individual).reshape(-1)
    prediction_scores = np.dot(features_dense, individual)
    predictions = prediction_scores > 0.5
    true_positives = np.sum((predictions == 1) & (labels == 1))
    true_negatives = np.sum((predictions == 0) & (labels == 0))
    false_positives = np.sum((predictions == 1) & (labels == 0))
    false_negatives = np.sum((predictions == 0) & (labels == 1))
    return true_positives * 2 + true_negatives - false_positives - 2 * false_negatives
```

4. **Selection:** Roulette wheel selection based on fitness scores.

Function to delect parents to generate a new generation:

```
def select(population, fitness):
    fitness_shifted = fitness - np.min(fitness) + 1e-3
    probability = fitness_shifted / np.sum(fitness_shifted)
    indices = np.random.choice(np.arange(population_size), size=population_size, p=probability)
    return population[indices]
```

- 5. Crossover: Experimented with single-point, multi-point, and uniform crossover.
 - Single-point:

```
def single_point_crossover(parent1, parent2):
    point = np.random.randint(1, len(parent1) - 1)
    child1 = np.concatenate((parent1[:point], parent2[point:]))
    child2 = np.concatenate((parent2[:point], parent1[point:]))
    return child1, child2
```

• Multi-point:

```
def multi_point_crossover(parent1, parent2):
    num_points = np.random.randint(1, len(parent1) - 1)
    points = sorted(np.random.choice(range(1, len(parent1)), num_points, replace=False))
    child1 = np.copy(parent1)
    child2 = np.copy(parent2)
    for i in range(0, len(points), 2):
        if i < len(points) - 1:
            child1[points[i]:points[i+1]], child2[points[i]:points[i+1]] = child2[points[i]:points[i+1]], child1[points[i]:points
            return child1, child2</pre>
```

• Uniform:

```
def uniform_crossover(parent1, parent2, prob=0.5):
    child1 = np.copy(parent1)
    child2 = np.copy(parent2)
    for i in range(len(parent1)):
        if np.random.rand() < prob:
            child1[i], child2[i] = child2[i], child1[i]
    return child1, child2</pre>
```

6. **Mutation**: Bit-flip mutation with mutation rates of 0.01, 0.05, and 0.1.

```
mutation rates = [0.01, 0.05, 0.1]
for mutation_rate in mutation_rates:
   print(f"Mutation Rate: {mutation rate}")
   best_chromosome = None
   best_fitness_score = -np.inf
   for generation in range(int(30)):
        fitness = np.array([calculate_fitness(ind, x_train_preprocessed, y_train) for ind in population])
       if np.max(fitness) > best_fitness_score:
           best_fitness_score = np.max(fitness)
           best_chromosome = population[np.argmax(fitness)]
        population = select(population, fitness)
        next_population = []
        for i in range(0, population_size, 2):
           parent1, parent2 = population[i], population[i+1]
            child1, child2 = crossover(parent1, parent2)
           next_population.extend([child1, child2])
        population = np.array([mutate(ind, mutation_rate) for ind in next_population])
        print(f"Generation {generation}: Best Fitness - {best_fitness_score}")
```

7. **Generations**: Iterated for a fixed number of generations (30).

3. INSIGHTS FROM THE CODE:

Fitness Function Choice:

The fitness function prioritizes accurate identification of both normal and intrusive network traffic instances by considering true positives and true negatives while penalizing false positives and false negatives.

Preprocessing Impact:

Standardization and one-hot encoding ensure that numerical and categorical features are appropriately scaled and represented for the genetic algorithm. This preprocessing step enhances the algorithm's ability to converge to optimal solutions.

Population Size Consideration:

Experimentation with different population sizes (50, 100, 200) revealed varying convergence speeds and solution quality. Larger population sizes may lead to faster convergence but require more computational resources.

Mutation Rate Impact:

Mutation rates of 0.01, 0.05, and 0.1 were tested, showing varied effects on convergence and solution quality. Lower mutation rates facilitate more stable convergence, while higher rates introduce more randomness, potentially aiding exploration.

Crossover Type Variation:

Experimentation with single-point, multi-point, and uniform crossover revealed differences in convergence behavior and solution quality. Each crossover type introduces different levels of exploration and exploitation, influencing algorithm convergence.

SUMMARY OF EVALUATIONS:

Evaluation Parameters	Value/Choice
Dataset	UNSW-NB15
Preprocessing	Standardization and One-Hot Encoding
Genetic Algorithm	Population Size: 50, 100, 200 Chromosome Length: Determined by features Fitness Function: Custom Selection: Roulette Wheel Crossover: Single-Point, Multi-Point, Uniform Mutation Rate: 0.01, 0.05, 0.1 Generations: 30
Insights from Code	 Fitness Function Choice Preprocessing Impact Population Size Consideration Mutation Rate Impact Crossover Type Variation

CONCLUSION:

The implementation of the genetic algorithm for intrusion detection on the UNSW-NB15 dataset demonstrates its effectiveness in optimizing feature selection. Various parameters such as population size, mutation rate, crossover type, and preprocessing techniques significantly impact the algorithm's convergence behavior and solution quality. Careful selection and fine-tuning of these parameters are crucial for achieving optimal performance. Further experimentation and

analysis may be required to explore additional optimization strategies and improve intrusion
detection accuracy