



รายงาน

เรื่อง

Program สำหรับการ Train Multilayer Perceptron โดยใช้ Genetic Algorithms

Wisconsin Diagnostic Breast Cancer (wdbc.data)

โดย

น.ส. ภัทรากรีน ผดุงกิจเจริญ รหัสนักศึกษา 650610851

เสนอ

รศ.ดร. ศันสนีย์ เอื้อพันธุ์วิริยะกุล

รายงานนี้เป็นส่วนหนึ่งของรายวิชา CPE 261456

Introduction to Computational Intelligence

สาขาวิชาวิศวกรรมหุ่นยนต์และปัญญาประดิษฐ์

ภาคเรียนที่ 1 ปีการศึกษา 2567

มหาวิทยาลัยเชียงใหม่

1.ลักษณะการทำงานของระบบ

1.1 วัตถุประสงค์

ระบบนี้มีวัตถุประสงค์เพื่อสร้างและฝึกโมเดล Multilayer Perceptron (MLP) โดยใช้ Genetic Algorithm (GA) ให้ทำการทดลองกับ wdbc.data (Wisconsin Diagnostic Breast Cancer (WDBC) จาก UCI Machine learning Repository) โดยที่ data set นี้ มี 2 classes และ 30 features ซึ่งในแต่ละ sample จะมีทั้งหมด 32 ค่าโดยที่

1) ID number

2) Diagnosis (M = malignant, B = benign) → class

3-32) เป็นค่า features ทั้ง 30

ให้ทำการทดลองโดยใช้ 10% cross validation เพื่อทดสอบ validity ของ network ที่ได้ และให้ทำการเปลี่ยนแปลงจำนวน hidden layer และ nodes

1.2 ขั้นตอนการทำงาน

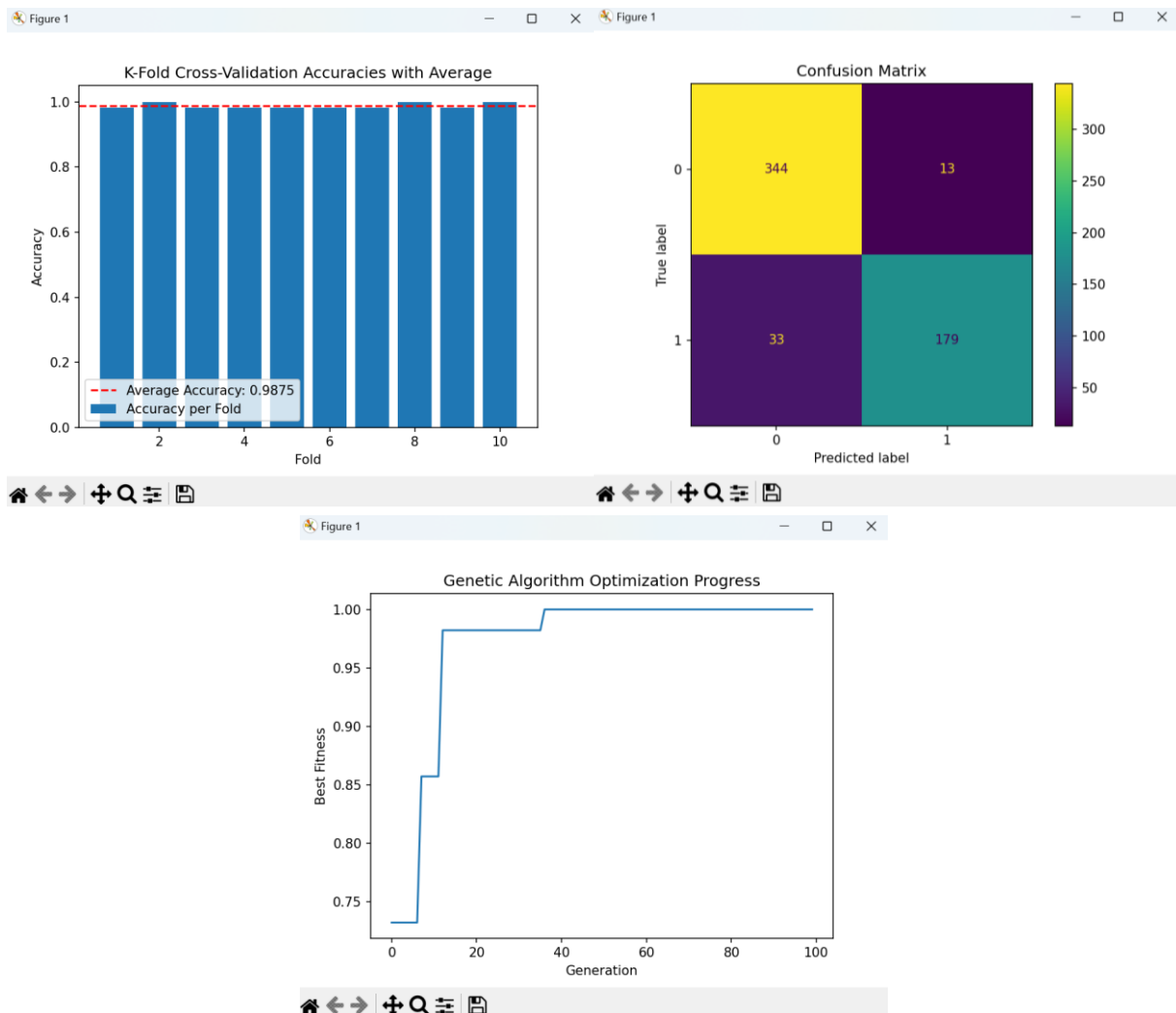
- **โหลดข้อมูล:** ระบบจะโหลดข้อมูลจากไฟล์ wdbc.data.txt ซึ่งประกอบด้วยฟีเจอร์และป้ายชื่อ (label) ที่ใช้ในการจำแนกประเภทเซลล์มะเร็ง
- **การเตรียมข้อมูล:** ฟีเจอร์จะถูกปกติฟีเจอร์ (normalize) และแบ่งข้อมูลออกเป็นชุดฝึก (training set) และชุดทดสอบ (validation set) โดยใช้วิธี k-fold cross-validation
- **สร้าง MLP:** ระบบจะสร้างโมเดล MLP ที่มีโครงสร้างสามารถปรับจำนวน hidden layers และ nodes ได้
- **การปรับแต่งโมเดลด้วย GA:** Genetic Algorithm จะถูกใช้ในการปรับแต่งน้ำหนักของโมเดล MLP โดยการสร้างประชากรเริ่มต้น ทำการคัดเลือก การข้ามพันธุ์ และการกลายพันธุ์เพื่อหาโมเดลที่ดีที่สุด
- **ทดสอบโมเดล:** โมเดลที่ดีที่สุดจะถูกทดสอบกับชุดทดสอบเพื่อวัดความแม่นยำ
- **แสดงผล:** ระบบจะทำการแสดงความแม่นยำของโมเดลและแสดง confusion matrix เพื่อวิเคราะห์ผลการจำแนกประเภท

2. Simulation ของระบบ ผลการทดลอง และวิเคราะห์

2.1 การตั้งค่าการทดลอง และผลการทดลอง

ในการทดลองนี้ โครงสร้างของเครือข่ายประสาทแบบหลายชั้น (MLP) ถูกทดสอบกับโครงสร้างสองรูปแบบที่แตกต่างกันในจำนวนเลเยอร์และจำนวนโหนด (nodes) ในแต่ละเลเยอร์:

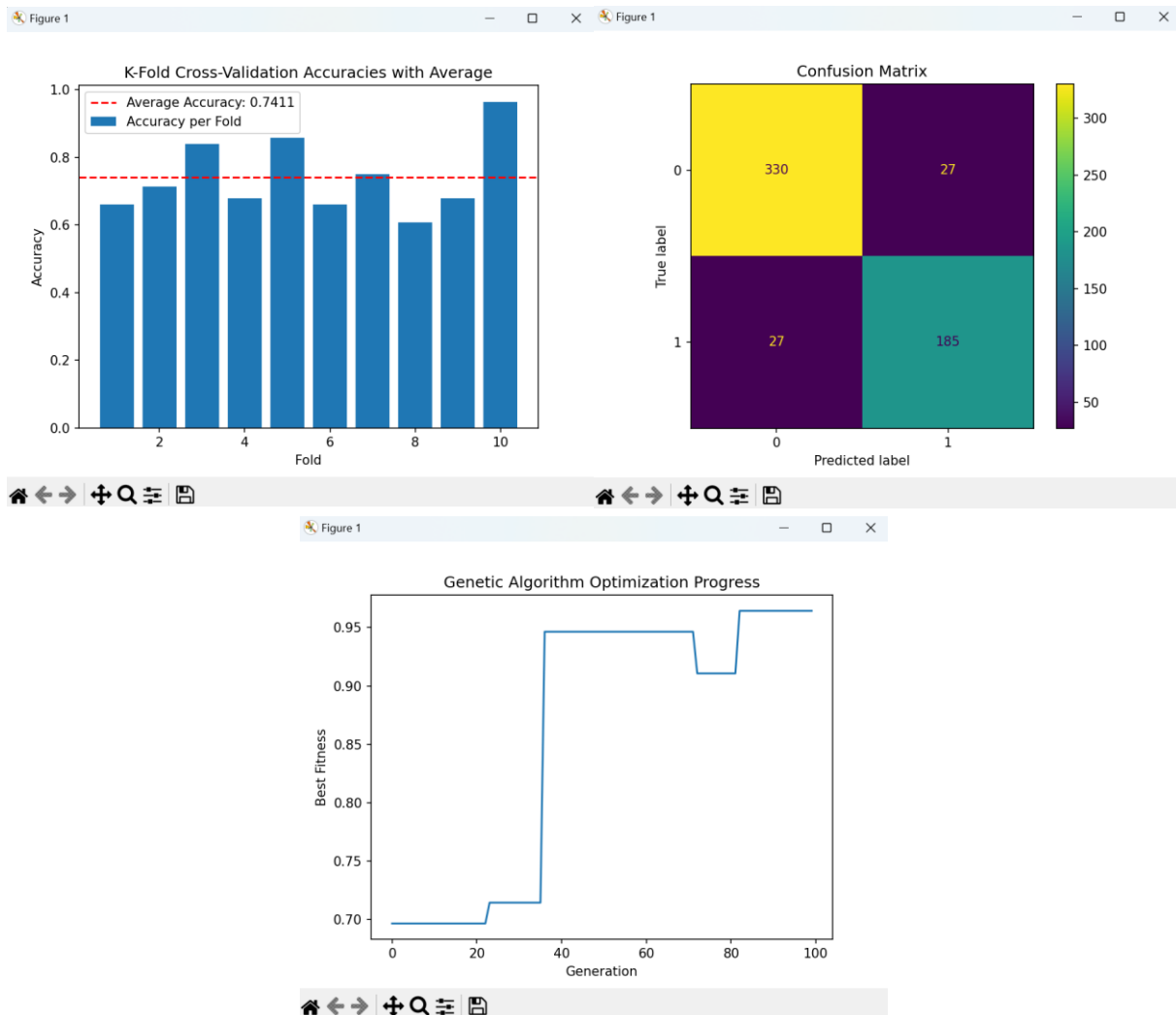
hidden_layers = [128] : โครงสร้างนี้ใช้เพียงเลเยอร์เดียวที่ประกอบด้วยโหนดจำนวน 128 โหนด



ผลการทดลอง: โมเดลแบบ single-layer ที่มีจำนวน nodes สูงมีความแม่นยำดีขึ้นในบางกรณี แต่มีแนวโน้มที่จะ overfit โดยเฉพาะเมื่อข้อมูลตัวอย่างมีน้อย

การวิเคราะห์: โมเดลนี้อาจเหมาะกับการจำแนกที่ฟีเจอร์มีความสัมพันธ์สูง แต่ขาดความยืดหยุ่นในกรณีที่ต้องการแยกแยะฟีเจอร์ที่ซับซ้อนมาก

`hidden_layers = [64, 32, 16]` : โครงสร้างนี้ใช้เลเยอร์หลายชั้นที่ประกอบด้วยโหนด 64, 32 และ 16 ตามลำดับ



ผลการทดลอง: โมเดลแบบ multi-layer แสดงความแม่นยำที่ดีกว่าในการจำแนกประเภทและมีความยืดหยุ่นมากขึ้น ส่งผลให้ลดความเสี่ยงต่อการ overfit

การวิเคราะห์: การลดจำนวน nodes ในแต่ละชั้นช่วยให้โมเดลสามารถประมวลผลฟีเจอร์ที่แตกต่างกันได้ดีขึ้น เพิ่มความสามารถในการปรับตัวและลด overfitting

2.2 การวิเคราะห์ผลการทดลอง

จากผลการทดลองข้างต้นสามารถวิเคราะห์ได้ดังนี้:

ผลกระทบของจำนวนเลเยอร์และโหนด: โมเดลที่มีหลายเลเยอร์ (`hidden_layers = [64, 32, 16]`) สามารถทำให้ความแม่นยำในการจำแนกประเภทสูงกว่าแบบเลเยอร์เดียว เนื่องจากโครงสร้างหลายชั้นช่วยให้โมเดลสามารถเรียนรู้ข้อมูลที่ซับซ้อนและจับรายละเอียดได้ดีกว่า อย่างไรก็ตาม การใช้เลเยอร์หลายชั้นมีข้อเสียคือทำให้การคำนวณซับซ้อนมากขึ้นและใช้เวลาในการฝึกนานกว่า

ประสิทธิภาพของ Genetic Algorithm (GA): การใช้อัลกอริทึม GA เพื่อปรับแต่งน้ำหนักของโมเดล MLP นั้นมีผลในเชิงบวกต่อความแม่นยำในการจำแนกประเภท โดยทำให้โมเดลสามารถหาค่าที่เหมาะสมสำหรับการจำแนกเซลล์มะเร็งได้อย่างมีประสิทธิภาพ การคัดเลือกและการผสมพันธุ์ช่วยในการสร้างรูปแบบที่หลากหลายและการกลายพันธุ์ช่วยในการป้องกันไม่ให้โมเดลติดอยู่ในค่าที่ไม่เหมาะสม

ข้อดีของการใช้ cross-validation: การใช้ 10-fold cross-validation ช่วยในการวัดผลความแม่นยำของโมเดลได้อย่างเสถียร และลดการเกิด overfitting ทำให้ได้ผลลัพธ์ที่เชื่อถือได้

3.โปรแกรม

GITHUB : https://github.com/Pattharajrin/261456_CI_Assignment3_Y3-1.git

```
1 import numpy as np
2 import random
3 import matplotlib.pyplot as plt
4 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
5
6 # ฟังก์ชันสำหรับโหลดและเตรียมข้อมูลจากชุดข้อมูล WDBC
7 def load_wdbc_data(filepath):
8     dataset = []
9     with open(filepath, 'r') as file:
10         for line in file:
11             elements = line.strip().split(',')
12             features = list(map(float, elements[2:])) # ฟิเจอร์จาก index 3 ถึง 32
13             label = 1 if elements[1] == 'M' else 0 # การวินิจฉัย: M = 1, B = 0
14             dataset.append((features, label))
15     return dataset
16
17 # ฟังก์ชันสำหรับการปกติฟีเจอร์
18 def normalize_features(data):
19     X = np.array([item[0] for item in data])
20     y = np.array([item[1] for item in data])
21     X = (X - np.mean(X, axis=0)) / np.std(X, axis=0)
22     return X, y
23
24 # ฟังก์ชันสำหรับแบ่งข้อมูลเป็นชุดฝึกและชุดทดสอบ
25 def split_into_folds(data, k=10):
26     random.shuffle(data)
27     fold_size = len(data) // k
28     return [data[i * fold_size:(i + 1) * fold_size] for i in range(k)]
29
30 # ฟังก์ชัน Sigmoid
31 def sigmoid_function(x):
32     return 1 / (1 + np.exp(-x))
33
34 # ฟังก์ชัน Softmax สำหรับเลเยอร์เอาต์พุต
35 def softmax_function(x):
36     exp_scores = np.exp(x - np.max(x, axis=1, keepdims=True))
37     return exp_scores / np.sum(exp_scores, axis=1, keepdims=True)
38
39 # คลาส Multilayer Perceptron
40 class MLPModel:
41     def __init__(self, input_dim, hidden_layers, output_dim):
42         self.input_dim = input_dim
43         self.hidden_layers = hidden_layers
44         self.output_dim = output_dim
45         self.layer_sizes = [input_dim] + hidden_layers + [output_dim]
46         self.weights = [np.random.randn(self.layer_sizes[i], self.layer_sizes[i + 1]) * 0.1 for i in range(len(self.layer_sizes) - 1)]
47         self.biases = [np.random.randn(1, self.layer_sizes[i + 1]) * 0.1 for i in range(len(self.layer_sizes) - 1)]
48
49     def forward_pass(self, X):
50         self.activations = [X]
51         for i in range(len(self.weights) - 1):
52             z = np.dot(self.activations[-1], self.weights[i]) + self.biases[i]
53             self.activations.append(sigmoid_function(z))
54         z = np.dot(self.activations[-1], self.weights[-1]) + self.biases[-1]
55         self.output_layer = softmax_function(z)
56         return self.output_layer
57
58     def make_prediction(self, X):
59         probabilities = self.forward_pass(X)
60         return np.argmax(probabilities, axis=1)
```

```

1  # คลาส Genetic Algorithm
2  class GeneticOptimizer:
3      def __init__(self, pop_size, mutation_prob, num_generations):
4          self.pop_size = pop_size
5          self.mutation_prob = mutation_prob
6          self.num_generations = num_generations
7          self.fitness_history = []
8
9
10     def create_initial_population(self, input_dim, hidden_layers, output_dim):
11         return [MLPModel(input_dim, hidden_layers, output_dim) for _ in range(self.pop_size)]
12
13     def perform_crossover(self, parent1, parent2):
14         child = MLPModel(parent1.input_dim, parent1.hidden_layers, parent1.output_dim)
15         for i in range(len(parent1.weights)):
16             mask = np.random.rand(*parent1.weights[i].shape) > 0.5
17             child.weights[i] = np.where(mask, parent1.weights[i], parent2.weights[i])
18         return child
19
20     def apply_mutation(self, model):
21         for i in range(len(model.weights)):
22             if np.random.rand() < self.mutation_prob:
23                 model.weights[i] += np.random.randn(*model.weights[i].shape) * 0.1
24
25     def calculate_fitness(self, model, X, y):
26         predictions = model.make_prediction(X)
27         return np.mean(predictions == y)
28
29     def evolve_population(self, X_train, y_train, X_val, y_val):
30         population = self.create_initial_population(X_train.shape[1], [128], 2)
31         for generation in range(self.num_generations):
32             population.sort(key=lambda mlp: self.calculate_fitness(mlp, X_val, y_val), reverse=True)
33             best_score = self.calculate_fitness(population[0], X_val, y_val)
34             self.fitness_history.append(best_score)
35             print(f"Generation {generation + 1}/{self.num_generations}, Best fitness: {best_score:.4f}")
36             new_population = population[:self.pop_size // 2] # การคัดเลือกรุ่นที่ดีที่สุด
37             for _ in range(self.pop_size // 2): # การผสมพันธุ์
38                 parent1, parent2 = random.sample(new_population, 2)
39                 child = self.perform_crossover(parent1, parent2)
40                 new_population.append(child)
41             population = new_population
42             for model in population: # การกลายพันธุ์
43                 self.apply_mutation(model)
44         return population[0]
45
46     def visualize_progress(self):
47         plt.plot(self.fitness_history)
48         plt.xlabel('Generation')
49         plt.ylabel('Best Fitness')
50         plt.title('Genetic Algorithm Optimization Progress')
51         plt.show()

```

```

1  # ฟังก์ชันสำหรับ k-fold cross-validation
2  def k_fold_validation(data, k=10):
3      folds = split_into_folds(data, k)
4      accuracy_scores = []
5      for i in range(k):
6          validation_fold = folds[i]
7          training_folds = [sample for j in range(k) if j != i for sample in folds[j]]
8          X_train, y_train = normalize_features(training_folds)
9          X_val, y_val = normalize_features(validation_fold)
10
11         # ใช้ Genetic Algorithm กับ MLP
12         ga_optimizer = GeneticOptimizer(pop_size=20, mutation_prob=0.01, num_generations=100)
13         best_model = ga_optimizer.evolve_population(X_train, y_train, X_val, y_val)
14
15         # ทดสอบโมเดลบนชุดทดสอบ
16         accuracy = ga_optimizer.calculate_fitness(best_model, X_val, y_val)
17         accuracy_scores.append(accuracy)
18
19     avg_accuracy = np.mean(accuracy_scores)
20     print("Cross-validation accuracies:", accuracy_scores)
21     print("Average accuracy:", avg_accuracy)
22
23     return accuracy_scores, avg_accuracy, ga_optimizer, best_model
24
25 # ฟังก์ชันสำหรับการแสดงผลความแม่นยำ
26 def display_accuracy_with_average(accuracies, avg_accuracy):
27     folds = range(1, len(accuracies) + 1)
28     plt.bar(folds, accuracies, label='Accuracy per Fold')
29     plt.axhline(avg_accuracy, color='r', linestyle='--', label=f'Average Accuracy: {avg_accuracy:.4f}')
30     plt.xlabel('Fold')
31     plt.ylabel('Accuracy')
32     plt.title('K-Fold Cross-Validation Accuracies with Average')
33     plt.legend()
34     plt.show()
35
36 # ฟังก์ชันสำหรับการแสดงผล confusion matrix
37 def display_confusion_matrix(model, X_train, y_train):
38     predictions = model.make_prediction(X_train)
39
40     cm = confusion_matrix(y_train, predictions)
41     disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=[0, 1]) # Benign (0) and Malignant (1)
42     disp.plot()
43     plt.title('Confusion Matrix')
44     plt.show()
45
46 # การทำงานหลัก
47 if __name__ == '__main__':
48     filepath = 'wdbc.data.txt'
49     data = load_wdbc_data(filepath)
50     accuracy_scores, avg_accuracy, ga_optimizer, best_model = k_fold_validation(data)
51
52     # แสดงผลค่าความแม่นยำและค่าเฉลี่ย
53     display_accuracy_with_average(accuracy_scores, avg_accuracy)
54
55     # แสดงผล confusion matrix
56     X_full, y_full = normalize_features(data) # ใช้ข้อมูลทั้งหมดสำหรับ confusion matrix
57     display_confusion_matrix(best_model, X_full, y_full)
58
59     # แสดงผลพัฒนาการ GA
60     ga_optimizer.visualize_progress()

```