

Chapter 10 Genetic Algorithm

Associate Professor Yachai Limpiyakorn, Ph.D.

Principles of Genetics

ในทางชีววิทยา เซลล์แต่ละเซลล์ในพืชชั้นสูงและสัตว์ประกอบด้วยนิวเคลียส (nucleus) 1 เซลล์ แต่ละนิวเคลียสประกอบด้วยโครโมโซม (chromosome) จำนวนหนึ่ง โดยโครโมโซมจะอยู่กันเป็นคู่มาจากพ่อและแม่อย่างละเส้น โครโมโซมแต่ละเส้นจะมียีน (gene) เป็นตัวกำหนดลักษณะถ่ายทอดทางพันธุกรรมของสิ่งมีชีวิต ในขณะที่มีการจับคู่กันของโครโมโซมอาจเกิด การไขว้เปลี่ยน (crossover) ซึ่งเป็นการที่ยีนจากโครโมโซมพ่อ แม่สลับเปลี่ยนกัน ทำให้เกิดโครโมโซมใหม่ขึ้น 2 คู่ และในขณะที่เซลล์แบ่งตัวจะเกิด กระบวนการ คัดลอกโครโมโซม (chromosome copying) ซึ่งบางครั้งจะมีการเปลี่ยนแปลงของยีนที่มาจากยีนพ่อแม่ เกิดเป็นยีนที่ไม่เคยมีมาก่อน เราเรียกการเกิดยีนลักษณะนี้ว่า การกลายพันธุ์ (mutation)

เราอาจกล่าวได้ว่า การคัดเลือกโดยธรรมชาติ (Natural Selection) เกิดจากการเปลี่ยนแปลงทางพันธุกรรมอันเป็นผลมาจาก การไขว้เปลี่ยน (crossover) ของ ลักษณะทางพันธุกรรม และ การกลายพันธุ์ (mutation)

2110773-10 2/2566

2

Natural Selection

- Charles Darwin ได้อธิบายการสืบทอดของสิ่งมีชีวิตด้วยกฎ **Evolution through Natural Selection** :-

- สิ่งมีชีวิตมีแนวโน้มที่จะสืบทอดลักษณะพิเศษให้ลูกหลาน
- ธรรมชาติจะผลิตสิ่งมีชีวิตที่มีลักษณะพิเศษแตกต่างไปจากเดิม
- สิ่งมีชีวิตที่เหมาะสมที่สุด (fittest) หรือที่มีลักษณะพิเศษที่ธรรมชาติพอใจมากที่สุด มีแนวโน้มที่จะมีลูกหลานมากกว่าตัวที่ไม่เหมาะสม ดังนั้น ประชากรจะโน้มเอียงไปทางตัวที่เหมาะสม
- การเปลี่ยนแปลงจะสะสมไปเรื่อยและเกิด species ใหม่ที่เหมาะสมกับสภาพแวดล้อม เมื่อเวลาผ่านไปนานๆ

2110773-10 2/2566

3

Genetic Algorithm (GA)

- Goldberg and Holland, 1988
- GA has its core idea from Darwin's theory of natural evolution "survival of the fittest"
- one of random-based *evolutionary algorithms* (EAs)
- search-based optimization technique **GA พัฒนาการ optimize**
- optimization → "how to find the best value for k that maximizes the performance of kNN classifier?"
- in order to find a solution, random changes applied to the current solutions to generate new ones.

2110773-10 2/2566

4

How GA works

- GA works on a population consisting of some solutions
- population size is the number of solutions
- each solution is called individual/ hypothesis
- each individual solution has a chromosome
- chromosome is represented as a set of parameters (features) that defines the individual
- each chromosome has a set of genes
- each gene is represented by somehow such as a string of 0s and 1s

2110773-10 2/2566

5

Hypothesis/ Individual Representation

Hypothesis: IF (Haircolor =black \vee red) \wedge (Eyecolor = dark) THEN Sunburn = negative

	Haircolor	Eyecolor	Sunburn
Bit String:	110	10	01

2 บิต = 2 bit

Hypothesis: IF Haircolor =blonde THEN Sunburn = positive

	Haircolor	Eyecolor	Sunburn
Bit String:	001	11	10

Note: Eyecolor = 11 \rightarrow Don't care

2110773-10 2/2566

6

Genetic Operators

◆ การไขว้เปลี่ยน (Crossover) เป็นการสร้างสายอักขระลูกหลาน 2 สาย จากสายอักขระพ่อแม่ 2 สาย โดยการคัดลอกบิตจากสายอักขระพ่อแม่ตามตำแหน่งที่กำหนดโดยหน้ากากไขว้เปลี่ยน (Crossover Mask) การไขว้เปลี่ยนทำได้หลายวิธีดังตัวอย่างที่แสดงข้างล่าง อาทิ การไขว้เปลี่ยน 1 ตำแหน่ง (Single-point Crossover) การไขว้เปลี่ยน 2 ตำแหน่ง (Two-point crossover) การไขว้เปลี่ยนยูนิฟอร์ม (Uniform Crossover) เพื่อให้ได้ความหลากหลาย (Diversity)

◆ การกลายพันธุ์ (Mutation) เป็นการสร้างการเปลี่ยนแปลงต่อสายอักขระลูกหลาน โดยสุ่มเลือกเปลี่ยนค่าบิตใดบิตหนึ่งดังแสดงในตัวอย่างข้างล่าง

2110773-10 2/2566

7

Parent Strings Crossover Mask Offsprings

Single-point crossover:

11101001000	11111000000	11101010101
00001010101		00001001000

Two-point crossover:

11101001000	00111110000	11001011000
00001010101		00101000101

Uniform crossover: *คือ ตำแหน่งการไขว้*

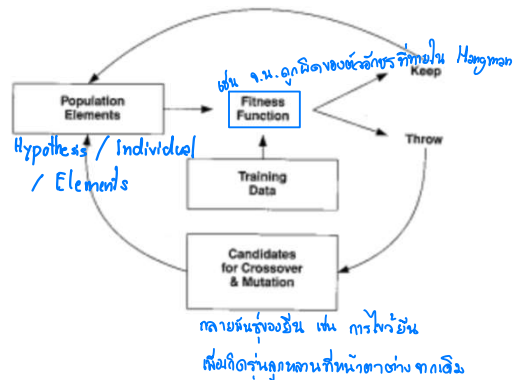
11101001000	10011010011	10001000100
00001010101		01101011001

Point mutation: *กลายพันธุ์ สลับค่าบิต*

11101001000		11101011000
-------------	--	-------------

2110773-10 2/2566

8



GA Overview

การเรียนรู้จำลอง
วิวัฒนาการของ
สิ่งมีชีวิต
(Biological
Evolution)

2110773-10 2/2566

9

GA Parameters

- Fitness function for ranking candidate patterns
- Stopping Criteria:
 - Maximum of hypotheses' fitness values \geq fitness threshold
 - Max fitness does not change after many generations
 - Reach fixed number of iterations of learning process
- Size of population p to be maintained
- Ratio of population r to be replaced at each generation
- Mutation rate m

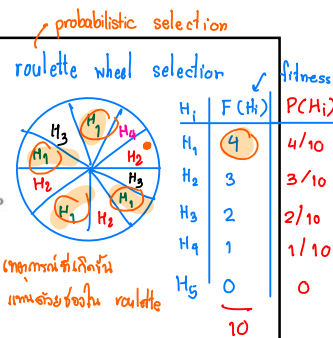
2110773-10 2/2566

10

GA (Fitness, Fitness_threshold, p, r, m)

- ◆ Initialize: $P \leftarrow p$ random hypotheses
- ◆ Evaluate: for each h in P compute $Fitness(h)$
- ◆ While $[max_i Fitness(h_i)] < Fitness_threshold$
 1. Select: Probabilistically select $(1-r)p$ individuals of P to add to P_s

$$Prob(h_j) = Fitness(h_j) / \sum_i Fitness(h_i) ; j=1, \dots, p$$
 2. Crossover: Probabilistically select $rp/2$ pairs of hypotheses from P .
For each pair $\langle h_1, h_2 \rangle$, produce two offspring by applying Crossover operator. Add all offspring to P_s .
 3. Mutate: Invert a randomly selected bit in mp random individuals of P_s .
 4. Update: $P \leftarrow P_s$.
 5. Evaluate: for each h in P compute $Fitness(h)$
- ◆ Return the hypothesis from P that has the highest fitness.



2110773-10 2/2566

11

12

Fitness Function and Selection

โดยทั่วไป โครโมโซมจะถูกเลือกแบบสุ่มเพื่อสร้างประชากรรุ่นใหม่ด้วยค่าความน่าจะเป็นตามสมการที่แสดงใน GA Algorithm ข้างต้น ซึ่งวิธีการเลือกดังกล่าว เรียกว่า Fitness Proportionate selection หรือ Roulette wheel selection เราจะสังเกตได้ว่า ฟังก์ชันค่าความเหมาะสม (fitness function) จะเป็นตัวกำหนดความน่าจะเป็นที่โครโมโซมจะอยู่รอดในประชากรรุ่น (generation) ถัดไป ซึ่งโครโมโซมที่มีค่าความเหมาะสมสูงจะมีโอกาสอยู่รอดมากกว่าโครโมโซมเส้นอื่นๆ แต่ก็ไม่ได้หมายความว่า โครโมโซมที่มีค่าความเหมาะสมสูงจะถูกเลือกทุกครั้ง ขึ้นอยู่กับการสุ่มค่า

2110773-10 2/2566

หน้าที่ยื่น เหลือเงิน ไม่ค่อยหลากหลาย

Crowding เป็นปรากฏการณ์ที่สมาชิกบางตัวในประชากร ซึ่งมีค่าความเหมาะสม ถูก reproduce อย่างรวดเร็ว ทำให้ความหลากหลาย (diversity) ของประชากรลดลง การแก้ปัญหา Crowding สามารถทำได้โดย

- เปลี่ยนวิธีการเลือกความเหมาะสม (Altering fitness selection) โดยใช้ Tournament Selection หรือ Rank Selection
- ใช้วิธีการ Fitness sharing หมายถึง ลดค่าความเหมาะสมของสมาชิกในประชากรรุ่นหนึ่งที่มีความคล้ายกัน (the fitness of a member is reduced by the presence of similar members) *ลดโอกาสที่จะถูกสุ่มเข้ามาใหม่*

2110773-10 2/2566

13

Tournament selection – เลือกสมาชิก 1 คู่จากประชากรรุ่นปัจจุบันแบบสุ่ม แล้วจึงสุ่มเลือกระหว่างสมาชิกคู่นั้น โดยสมาชิกตัวที่มีค่าความเหมาะสมสูงกว่าจะถูกสุ่มเลือกด้วยความน่าจะเป็นที่กำหนดไว้ล่วงหน้า p ส่วนสมาชิกตัวที่มีค่าความเหมาะสมต่ำกว่าจะถูกสุ่มเลือกด้วยความน่าจะเป็น $(1-p)$ *คนที่ถูกเลือกขึ้น p ไม่ถูกเลือกขึ้น $1-p$ ซึ่ง p เหนือคือ > 0.5*

Rank selection – เป็นวิธีที่ใช้ควบคุมการเลือกโครโมโซมโดยไม่สนใจค่าความเหมาะสมของโครโมโซมว่ามีค่าเท่าไร แต่จะใช้ค่าความเหมาะสมเพียงแต่จัดลำดับเรียงโครโมโซมตามค่าความเหมาะสมที่มีค่าสูงที่สุดจนถึงต่ำสุด จากนั้น กำหนดค่าคงที่ p เป็นความน่าจะเป็นที่โครโมโซมลำดับที่ 1 จะถูกเลือก และเป็นความน่าจะเป็นที่โครโมโซมลำดับที่ 2 จะถูกเลือกเมื่อโครโมโซมลำดับที่ 1 ไม่ถูกเลือก และเป็นความน่าจะเป็นที่โครโมโซมลำดับที่ 3 จะถูกเลือกเมื่อลำดับที่ 1 และ 2 ไม่ถูกเลือก เรียงไปจนกระทั่งถึงลำดับสุดท้ายซึ่งจะถูกเลือกเมื่อลำดับก่อนหน้าไม่ถูกเลือกเลย ดังนั้น ความน่าจะเป็นที่โครโมโซมลำดับที่ r จะถูกเลือก เท่ากับ $p(1-p)^{r-1}$; $r = 1, 2, 3, \dots$

2110773-10 2/2566

14

Rank Selection

Hypothesis	Fitness(h_i)	$P(h_i)$	Rank	Probability of Rank(h_i)
h_1	4	0.4	1	$=p = 0.667$
h_2	3	0.3	2	$=p(1-p) = 0.667 \times 0.333 = 0.222$
h_3	2	0.2	3	$=p(1-p)^2 = 0.667(0.333)^2 = 0.074$ <i>Ranking - 1</i>
h_4	1	0.1	4	$=p(1-p)^3 = 0.667(0.333)^3 = 0.025$
h_5	0	0	5	$=p(1-p)^4 = 0.667(0.333)^4 = 0.012$

2110773-10 2/2566

นี่คือให้ Hypothesis น้อยๆ ถูกเลือก

15

(No) Credit + Yes *เพราะ No* *No*

$H_2 = \langle 20-30k, \text{Yes}, \text{Male}, 30-39 \rangle$; $F(H_2) = \frac{2+0+2+0}{0+1+1+2+1} = \frac{4}{5} < 1$; *ไม่ผ่าน crossover*

Example GA Application *fitness $F(H_i) = \frac{N}{M}$*

ตัวอย่างชุดข้อมูลสอนการจำแนกประเภทลักษณะลูกค้าที่สนใจ Life Insurance Promotion

Training Instance	Income Range	Life Insurance Promotion	Credit Card Insurance	Sex	Age
1	30-40K	Yes	Yes	Male	30-39
2	30-40K	Yes	No	Female	40-49
3	50-60K	Yes	No	Female	30-39
4	20-30K	No	No	Female	50-59
5	20-30K	No	No	Male	20-29
6	30-40K	No	No	Male	40-49

2110773-10 2/2566

16

$H_3 = \langle ?, \text{No}, \text{No}, \text{Male}, 40-49 \rangle$; $F(H_3) = \frac{0+3+2+1}{0+2+1+1+1} = \frac{6}{5} > 1$ *ไม่ผ่าน crossover*

Fitness function

$$F(E_i) = N / (M+1)$$

where

N คือ จำนวนตัวอย่างในชุดข้อมูลสอนซึ่งอยู่ในคลาสเดียวกับ E_i ที่มีค่าคุณลักษณะตรงกัน

M คือ จำนวนตัวอย่างในชุดข้อมูลสอนซึ่งไม่อยู่ในคลาสเดียวกับ E_i ที่มีค่าคุณลักษณะตรงกัน

หมายเหตุ เพื่อป้องกันการหารด้วยค่าศูนย์ จึงบวกหนึ่งที่ตัวหารในฟังก์ชันความเหมาะสม

กฎการคัดเลือกสมาชิกของประชากรรุ่นถัดไป

- สมาชิกที่มีค่าความเหมาะสมน้อยกว่า ค่าขีดแบ่ง (threshold) ที่กำหนดไว้มีค่าเท่ากับ 1 จะถูกนำไป crossover หรือ mutation
- ถ้าสมาชิกทุกตัวมีค่าความเหมาะสมไม่น้อยกว่าค่าขีดแบ่ง จะสุ่มเลือกสมาชิกเพื่อไป crossover หรือ mutation
- ในประชากรทุกรุ่น จะต้องมีความเหมาะสมที่อยู่ในคลาส "Yes" และ "No" อย่างละ 2 ตัว

กำหนดขนาดประชากรแต่ละรุ่นเท่ากับ 4 ตัว

Income Range = ? → don't care

สมาชิกประชากรรุ่นแรก

An Initial Population for Supervised Genetic Learning

Population Element	Income Range	Life Insurance Promotion	Credit Card Insurance	Sex	Age
1	20-30K	No	Yes	Male	30-39
2	30-40K	Yes	No	Female	50-59
3	?	No	No	Male	40-49
4	30-40K	Yes	Yes	Male	40-49

ตัวอย่างการคำนวณ $F(E_1)$

- Income Range = 20-30K ตรงกับตัวอย่างสอนที่ 4 และ 5
- Credit Card Insurance = Yes ไม่ตรงกับตัวอย่างสอนใดๆ
- Sex = Male ตรงกับตัวอย่างสอนที่ 5 และ 6
- Age = 30-39 ไม่ตรงกับตัวอย่างสอนใดๆ
- จะได้ว่า ค่า $N = 2+0+2+0 = 4$
- Income Range = 20-30K ไม่ตรงกับตัวอย่างสอนใดๆ
- Credit Card Insurance = Yes ตรงกับตัวอย่างสอนที่ 1
- Sex = Male ตรงกับตัวอย่างสอนที่ 1
- Age = 30-39 ตรงกับตัวอย่างสอนที่ 1 และ 3
- จะได้ว่า ค่า $M = 0+1+1+2 = 4$

การไขว้เปลี่ยนเพื่อสร้างสมาชิกลักษณะใหม่

- $F(E_1) = 4/5 = 0.80$, $F(E_2) = 6/7 = 0.86$; $F(E_1)$ and $F(E_2) < 1$ → crossover
- $F(E_3) = 6/5 = 1.20$, $F(E_4) = 5/5 = 1.00$; $F(E_3)$ and $F(E_4) \geq 1$ → keep

crossover ยังไม่จบถึง class label

Population Element	Income Range	Life Insurance Promotion	Credit Card Insurance	Sex	Age
#1	20-30K	No	Yes	Male	30-39
#2	30-40K	Yes	No	Fem	50-59

Population Element	Income Range	Life Insurance Promotion	Credit Card Insurance	Sex	Age
#2	30-40K	Yes	Yes	Male	30-39
#1	20-30K	No	No	Fem	50-59