

การนำการแจกแจงความถี่ของข้อมูลทางสถิติ
แบบ Beta Distribution มาประยุกต์ใช้เพื่อหา
ค่าความเหมาะสมภายใต้ความแปรปรวน ความเบ้
และความโด่งที่แตกต่างกันสำหรับตัวแปรเชิงกายภาพ
ร่วมกับค่าสหสัมพันธ์ด้วยโปรแกรม MATLAB และการเปรียบ
เทียบกับการใช้ Gaussian Function

The Calculation of Suitability Score under Beta
Distribution Probability Density Function Curve using
Variance, Skewness, and Excess Kurtosis as Arguments
for Physical Quantities with MATLAB Software and the
Comparison with Gaussian Function

วิวรรษธร ฐิตศิริวิทย์

31 สิงหาคม 2563

สารบัญ

1	บทนำ (Introduction)	1
1.1	การแจกแจงความถี่ทางสถิติรูปแบบต่าง ๆ	1
2	พารามิเตอร์ของฟังก์ชัน (Function Parameters)	1
2.1	สัญลักษณ์ที่ใช้เขียนฟังก์ชันที่มีตัวแปรและพารามิเตอร์ควบคุม	1
3	นิยามของฟังก์ชันที่เกี่ยวข้องในทางสถิติ	1
3.1	ค่ากลางของข้อมูล	1
3.1.1	ค่าเฉลี่ยเลขคณิต (Mean, μ หรือ \bar{x})	1
3.1.2	ค่ามัธยฐาน (Median, Me)	2
3.1.3	ค่าฐานนิยม (Mode, Mo)	2
3.2	ส่วนเบี่ยงเบนมาตรฐานและความแปรปรวน	2
3.2.1	ส่วนเบี่ยงเบนมาตรฐาน (Standard Deviation, σ หรือ s)	2
3.2.2	ความแปรปรวน (Standard Deviation, σ^2 หรือ s^2)	2
3.3	ค่าความเบ้ (Skewness)	2
3.4	ค่าความโด่ง (Excess Kurtosis)	2
3.5	ฟังก์ชันแกมมา (Gamma Function, $\Gamma(x)$)	3
4	ฟังก์ชันการกระจายแบบสมมาตร (Symmetric Distribution)	3
4.1	Normal Distribution	3
4.2	Lorentzian Distribution (Cauchy Distribution/Breit-Wigner Distribution)	3
5	ฟังก์ชันการกระจายแบบอสมมาตร (Asymmetric Distribution)	4
5.1	Chi-Square (χ^2) Distribution	4
5.2	Gamma (γ) Distribution (Erlang Distribution)	4
5.3	Beta (β) Distribution	4
6	การใช้โปรแกรม MATLAB วิเคราะห์ข้อมูลเบื้องต้น	5
6.1	ชนิดของข้อมูล	5
6.2	การสร้างฟังก์ชันคำนวณ	5
7	การนำ Normal Distribution มาประยุกต์ใช้งาน	6
7.1	ที่มาของการประยุกต์และ การสร้าง Gaussian Function เพื่อวิเคราะห์ข้อมูล	6
7.2	ข้อได้เปรียบและข้อจำกัดของ Gaussian Function ใน Symmetric Distribution	6
7.3	การสร้างฟังก์ชัน Gaussian Suitability Score, S_N	7
8	การนำ Beta Distribution มาประยุกต์ใช้งาน	8
8.1	ที่มาของฟังก์ชันหาค่าความเหมาะสม S_β	8
8.2	โดเมนและเรนจ์ของฟังก์ชัน	9
8.2.1	โดเมนของตัวแปรและพารามิเตอร์คงที่	9
8.2.2	เรนจ์ของ S_β	10
8.3	การสร้างฟังก์ชัน Beta Distribution Suitability Score	10

9	เนื้อหาเพิ่มเติม	11
9.1	สัญลักษณ์ Knuth's up-arrow Notation	11
9.2	ตัวอย่างการสร้างฟังก์ชันคำนวณด้วย MATLAB (เก่า)	11

1 บทนำ (Introduction)

1.1 การแจกแจงความถี่ทางสถิติรูปแบบต่าง ๆ

การแจกแจงทางสถิติ (Statistic distribution) มีหลากหลายรูปแบบ หนึ่งในเกณฑ์การจำแนก คือ

1. การแจกแจงแบบสมมาตร (Symmetric distribution) เช่น Normal Distribution, Lorentzian Distribution ฯลฯ
คือ กราฟของชุดข้อมูลแบบฟังก์ชันแบบ Probability Density Function (PDF) ที่กระจายตัวอย่างสมมาตรในแนว $x = x_0$
2. การแจกแจงแบบไม่สมมาตร (Asymmetric distribution) เช่น Beta Distribution, Gamma Distribution, Chi-Square Distribution ฯลฯ
คือ กราฟของชุดข้อมูลแบบฟังก์ชัน PDF ที่มีค่าความเบ้ (Skewness, Skew) ในทางซ้าย/ขวา หรือ ลบ/บวก โดย x_0 จะไม่อยู่ในแนวสมมาตร

และอาจแบ่งได้เป็นแบบตามชุดข้อมูลแบบไม่ต่อเนื่อง (Discrete distribution) และแบบต่อเนื่อง (Continuous distribution) แต่ในที่นี้จะเน้นไปในหัวข้อ การแจกแจงแบบต่อเนื่อง ซึ่งอยู่ในรูปของ PDF และจะไม่กล่าวถึง Cumulative Probability Function (CDF)

2 พารามิเตอร์ของฟังก์ชัน (Function Parameters)

2.1 สัญลักษณ์ที่ใช้เขียนฟังก์ชันที่มีตัวแปรและพารามิเตอร์ควบคุม

สัญลักษณ์ของฟังก์ชัน ใช้รูป $y = f(x_1, x_2, x_3, \dots, x_n; a_1, a_2, a_3, \dots, a_n)$ โดยสิ่งที่อยู่ใน $f(\dots)$ นั้นคือ arguments, x_n คือ variable, a_n คือ fixed parameter ที่มีอยู่แล้ว ไม่ใช่หาความสัมพันธ์ เป็นตัวแปรควบคุมในการหา $x_n \mapsto f(x_n)$

3 นิยามของฟังก์ชันที่เกี่ยวข้องในทางสถิติ

3.1 ค่ากลางของข้อมูล

3.1.1 ค่าเฉลี่ยเลขคณิต (Mean, μ หรือ \bar{x})

ค่าเฉลี่ยเลขคณิต คือ ผลรวมของข้อมูล หารด้วยจำนวนชุดข้อมูล อาจด้วยน้ำหนักที่เท่าหรือไม่เท่า ขึ้นอยู่กับการเลือกใช้ค่าเฉลี่ยให้เหมาะสมกับชุด ว่าต้องการให้น้ำหนักกับข้อมูลใดเป็นพิเศษหรือไม่ ดังสมการที่ 1 และ 2 ตามลำดับ

$$\begin{aligned}\mu_x &= \frac{1}{N} \sum_{i=1}^N x_i \\ &= \frac{1}{\beta - \alpha} \int_{\alpha}^{\beta} f(x_i) dx\end{aligned}\tag{1}$$

$$\mu_{x,w} = \frac{1}{\sum w_i} \sum_{i=1}^N w_i x_i\tag{2}$$

ส่วนในโปรแกรม MATLAB ใช้คำสั่ง `mean(A)`

3.1.2 ค่ามัธยฐาน (Median, Me)

ค่ามัธยฐาน คือ ค่าที่อยู่ตรงกลางของช่วงข้อมูล แจกแจงตามความถี่ในช่วงต่อเนื่อง¹

3.1.3 ค่าฐานนิยม (Mode, Mo)

ค่าฐานนิยม คือ ค่าที่มีความถี่สูงสุดของช่วงข้อมูลต่อเนื่อง²

3.2 ส่วนเบี่ยงเบนมาตรฐานและความแปรปรวน

3.2.1 ส่วนเบี่ยงเบนมาตรฐาน (Standard Deviation, σ หรือ s)

ส่วนเบี่ยงเบนมาตรฐาน คือ ค่าที่ใช้การกระจายของข้อมูลจากค่าเฉลี่ย หากจากรากที่สองของความแปรปรวน คำนวณได้จาก

สมการที่ 3 ซึ่งในโปรแกรม MATLAB สามารถหาได้จาก `std(A)` โดย A เป็นเวกเตอร์ของข้อมูลประเภทหนึ่ง

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x - \mu)^2} \quad (3)$$

3.2.2 ความแปรปรวน (Standard Deviation, σ^2 หรือ s^2)

ความแปรปรวน คือ ตัวแปรที่ใช้วัดการกระจายของข้อมูล คิดจากการหาค่าเฉลี่ยของความต่างจากค่าเฉลี่ยกำลังสอง คำนวณได้จากสมการที่ 4 โดยในโปรแกรม MATLAB สามารถหาได้จาก `var(B)` โดย B เป็นเวกเตอร์ของข้อมูลประเภทหนึ่ง

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x - \mu)^2 \quad (4)$$

3.3 ค่าความเบ้ (Skewness)

ค่าความเบ้ ใช้บอกว่าข้อมูลส่วนใหญ่ (ชุดหรือช่วงที่มีความถี่สูง) ไปอยู่ที่ Interval ไตของกราฟการกระจาย

3.4 ค่าความโด่ง (Excess Kurtosis)

ค่า Excess Kurtosis ใช้บอกว่าชุดข้อมูลต่อเนื่องนั้น มีการกระจาย หรือการเป็นกลุ่มมากแค่ไหน เมื่อเทียบกับ Normal Distribution

¹ในการคำนวณมัธยฐานของช่วงข้อมูลที่ไม่ต่อเนื่อง ทั้งแบบแจกแจงความถี่และไม่แจกแจงความถี่ ต่างกับการหามัธยฐานของชุดข้อมูลแบบต่อเนื่อง (PDF) ขึ้นอยู่กับรูปแบบของการกระจายของข้อมูลว่าเป็นลักษณะใด

²ฐานนิยมแบบไม่ต่อเนื่องกับต่อเนื่อง ต่างกันในการทำงานเกี่ยวกับการหามัธยฐาน

3.5 ฟังก์ชันแกมมา (Gamma Function, $\Gamma(x)$)

ฟังก์ชันแกมมา คือ ฟังก์ชันที่ให้นิยามของแฟกทอเรียลที่นอกเหนือจากสมการที่ 5

$$n! = n(n-1)(n-2) \dots (3)(2)(1) ; n \in \mathbb{Z}^+ \quad (5)$$

ให้นิยามและสัญลักษณ์ดังสมการที่ 6

$$\begin{aligned} \Gamma(\xi) &= (\xi-1)! \\ &= \int_0^\infty e^{-t} t^{\xi-1} dt \end{aligned} \quad (6)$$

4 ฟังก์ชันการกระจายแบบสมมาตร (Symmetric Distribution)

4.1 Normal Distribution

Normal Distribution คือ การกระจายของข้อมูลรูปแบบที่ทั่วไปที่สุด (ที่มีการสุ่มอย่างแท้จริง, true random) สามารถใช้ได้หลากหลายกรณี กล่าวคือเป็นการแจกแจงความถี่ของข้อมูลในอุดมคติ โดยคำนึงถึงเพียงแค่ความแปรปรวน (σ^2) และค่าเฉลี่ยเลขคณิต (μ) ของข้อมูลเท่านั้น ทำให้ง่ายต่อการศึกษาและใช้งานในระดับเบื้องต้น โดยมี PDF ดังสมการที่ 7

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} ; \sigma > 0 \quad (7)$$

โดย Skew = 0, $\bar{x} = \mu$, $s^2 = \sigma^2$ และ *ExcessKurtosis* = 0

4.2 Lorentzian Distribution (Cauchy Distribution/Breit-Wigner Distribution)

Lorentzian Distribution คือ การกระจายข้อมูลซึ่งขึ้นกับตัวแปร Scale (γ) และตัวแปร ตำแหน่งของฐานนิยม (m หรือ x_0) สามารถปรับการกระจายของข้อมูลได้ มาจาก Normal Distribution ที่แตกต่างกัน 2 แบบที่ไม่ขึ้นกับกันและกัน ขณะที่ $\bar{x} = 0$ และ $s^2 = 1$ โดยมี PDF ดังสมการที่ 8

$$f(x; x_0, \gamma) = \frac{1}{\pi\gamma \left(1 + \left(\frac{x-x_0}{\gamma}\right)^2\right)} ; \gamma > 0 \quad (8)$$

โดย Skew หาค่าไม่ได้, $\bar{x} = 0$, s^2 หาค่าไม่ได้ และ *ExcessKurtosis* หาค่าไม่ได้

5 ฟังก์ชันการกระจายแบบอสมมาตร (Asymmetric Distribution)

5.1 Chi-Square (χ^2) Distribution

Chi-Square Distribution คือ การกระจายข้อมูลที่มักใช้ในการทดสอบสมมติฐาน มีความคล้ายคลึงกับ Normal Distribution และ Gamma Distribution ซึ่งขึ้นกับค่า chi-square (χ^2) และค่าองศาเสรี (Degree of freedom, k) โดยมี PDF ดังสมการที่ 9

$$f(x; k) = \frac{2^{(-\frac{k}{2})}}{\Gamma(\frac{k}{2})} x^{\frac{k}{2}-1} e^{-\frac{x}{2}} ; x > 0 \quad (9)$$

โดย $Skew = 2\sqrt{\frac{2}{k}}$, $\bar{x} = k$, $s^2 = 2k$ และ $ExcessKurtosis = \frac{12}{k}$

5.2 Gamma (γ) Distribution (Erlang Distribution)

Gamma Distribution มักถูกใช้ในการวิเคราะห์ปริมาณทางฟิสิกส์ที่มีการแจกแจงความถี่ ทั้งการศึกษา และการสรุปผลการทดลองที่ไม่เป็นแบบสุ่ม (non-random events) เช่น สถิติการรอคิวในช่วงเวลาใดเวลาหนึ่ง โดยมีความคล้ายคลึงกับการแจกแจงความถี่รูปแบบอื่น ๆ ด้วย และขึ้นกับตัวแปร α และ β เท่านั้น โดยมี PDF ดังสมการที่ 10

$$f(x; \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta} \left(\frac{x}{\beta}\right)^{\alpha-1} e^{-\frac{x}{\beta}} ; \alpha, \beta > 0 \quad (10)$$

โดย $Skew = \frac{2}{\sqrt{\alpha}}$, $\bar{x} = \alpha\beta$, $s^2 = \alpha\beta^2$ และ $ExcessKurtosis = \frac{6}{\alpha}$

5.3 Beta (β) Distribution

Beta Distribution คือการแจกแจงความถี่ที่ขึ้นกับตัวแปร α และ β เท่านั้น ซึ่งใช้บอกความเบ้และความโด่งได้ เหมือนกับ Gamma Distribution และมีลักษณะคล้ายคลึงกับ Bernoulli Distribution (Yes/No) ซึ่งเป็นแบบ Discrete Distribution

การแจกแจงความถี่ในลักษณะนี้มีความยืดหยุ่นต่อลักษณะตัวแปร และข้อมูลหลากหลายประเภท จึงนำมาประยุกต์ใช้ในงานวิเคราะห์ขั้นสูงได้มาก มีการใช้งานในวงกว้างเหมือนกับแบบ Normal Distribution

อีกลักษณะที่เด่น คือการแจกแจงข้อมูลแบบนี้ สามารถกำหนดช่วงของข้อมูลที่เก็บมาได้ให้เป็นช่วง และจำกัด Scale (γ) ให้เป็นค่าที่ต้องการได้ โดย Beta Distribution มี PDF ดังสมการที่ 11

$$f(x; \alpha, \beta) = \frac{1}{\left[\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}\right]} x^{\alpha-1} (1-x)^{\beta-1} ; \alpha, \beta > 0 \quad (11)$$

โดย $Skew = \frac{2(\beta-\alpha)\sqrt{\alpha+\beta+1}}{(\alpha+\beta+2)\sqrt{\alpha\beta}}$, $\bar{x} = \frac{\alpha}{\alpha+\beta}$, $s^2 = \frac{\alpha\beta}{(\alpha+\beta)^2(1+\alpha+\beta)}$
และ $ExcessKurtosis = \frac{3(\alpha+\beta+1)[\alpha\beta(\alpha+\beta-6)+2(\alpha+\beta)^2]}{\alpha\beta(\alpha+\beta+2)(\alpha+\beta+3)} - 3$

ในกรณีที่ $0 < \alpha, \beta < 1$ ทำให้ข้อมูลกระจายไปทางซ้ายและขวามากกว่าตรงกลาง แต่ในกรณีที่ $\alpha, \beta > 1$ จะทำให้ข้อมูลเป็นรูปโค้งคว่ำ

6 การใช้โปรแกรม MATLAB วิเคราะห์ข้อมูลเบื้องต้น

6.1 ชนิดของข้อมูล

ในโปรแกรม MATLAB ตัวแปรที่ใช้เป็นหลักจะอยู่ในรูปของอาร์เรย์ (Array) ซึ่งส่วนใหญ่ มักใช้เพียง 2 รูปแบบ คือ อาร์เรย์ 1 มิติ (เวกเตอร์หลัก/column vector) และอาร์เรย์ 2 มิติ (เมทริกซ์/matrix)

การเก็บข้อมูลเชิงปริมาณ (Quantitative Data) หรือในที่นี้ ข้อมูลเชิงตัวเลข (Numerical Data) จะอยู่ในรูปของค่าที่ถูกล้อมด้วยเครื่องหมายจุลภาค (Comma-separated value หรือ Comma-delimited value, CSV Format) คือเป็นลักษณะของหลัก (column) และเมื่อเก็บข้อมูลซ้ำ ๆ จะทบไปในทางแถว (row) ขณะที่จำนวนหลักคงที่ตลอดการเก็บ 1 ช่วง เช่น

```
253,25.49,65.36,1012.17,25.3
254,26.50,68.77,1012.95,24.1
255,26.55,69.01,1013.00,20.4
256,26.86,69.68,1013.24,15.2
```

ซึ่งเป็นตัวอย่างชุดข้อมูลที่เก็บมา³ ในรูปแบบ CSV (สกุลไฟล์ที่เก็บข้อมูลประเภทนี้ไม่มีความเกี่ยวข้องกับรูปแบบการจัดเรียงข้างใน หากไฟล์ถูกสร้างขึ้นจากไฟล์เปล่า) สามารถเทียบเคียงได้กับอาร์เรย์:

	col1	col2	col3	col4	col5
row1	253	25.49	65.36	1012.17	25.3
row2	254	26.50	68.77	1012.95	24.1
row3	255	26.55	69.01	1013.00	20.4
row4	256	26.86	69.68	1013.24	15.2

Table 1: Array A

อาร์เรย์นี้มีขนาด 4 แถว x 5 หลัก ซึ่งในโปรแกรม MATLAB สามารถหาขนาดของอาร์เรย์ได้ผ่านคำสั่ง `size(A)` ซึ่งค่าที่ออกมาจะเป็นในรูปแบบของเวกเตอร์แถว (row vector) เช่นในกรณีนี้ อาร์เรย์มีขนาด `[4,5]`

การบอกตำแหน่งของค่า เรียกว่า index จะบอกเป็นพิกัด (แถว,หลัก) ของอาร์เรย์ สามารถระบุเป็นจุดได้ `A(m,n)`, เลือกเวกเตอร์แถวทั้งช่วง `A(m,:)` หรือบางช่วง `A(m,a:b)` และเลือกเวกเตอร์หลักทั้งช่วง `A(:,n)` หรือบางช่วง `A(a:b,n)` หรือแม้แต่ดึงอาร์เรย์มาบางส่วน `A(a:b,c:d)` ก็ได้ ซึ่งส่วนใหญ่ในการวิเคราะห์ มักดึง column vector ออกจาก matrix เพื่อวิเคราะห์แยกทีละปัจจัยมากกว่า

6.2 การสร้างฟังก์ชันคำนวณ

ให้ไปศึกษาเอาเอง (หัวข้อ `function out = calc(path arguments)`) รวมถึงการดึง file path, การทำ composite function การสร้างสมการคำนวณ การวิเคราะห์แบบเส้น-เส้น การทำ fitting และการสร้าง loop

³การเขียนไฟล์ในลักษณะนี้ ใช้การ `rw-` แบบ `++a` (append if existed) ได้เลย ทำให้ข้อมูลต่อท้าย และเมื่อ algorithm ในการตรวจหาไฟล์ซ้ำที่มีอยู่แล้วเพื่อป้องกันการเขียนทับด้วย

7 การนำ Normal Distribution มาประยุกต์ใช้งาน

7.1 ที่มาของการประยุกต์และ การสร้าง Gaussian Function เพื่อวิเคราะห์ข้อมูล

ฟังก์ชันการแจกแจงความถี่โค้งแบบ Normal Distribution มีความเป็นอุดมคติในทางข้อมูลและรูปแบบการกระจายตัวที่สุด ทำให้ง่ายต่อการใช้งาน และการนำมาแก้ไข ดัดแปลงให้เป็น Non-normal distribution แบบอื่น ๆ ขณะที่ยังคงความเบ้เป็นศูนย์เหมือนเดิม เดิมทีฟังก์ชัน Normal Distribution จะอยู่ในรูปของสมการที่ 12

$$N(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (12)$$

ในการดัดแปลง เราอาจนำค่าคงที่ $\frac{1}{\sqrt{2\pi}\sigma}$, ตัวคูณ $\frac{1}{2\sigma^2}$ และ X-control variable $-\mu$ ออก จะได้เป็นสมการที่ 13 ซึ่งอยู่ในรูปมาตรฐานของฟังก์ชัน Gaussian Function (สามารถหาปริพันธ์จำกัดเขตตั้งแต่ $-\infty$ ถึง $+\infty$ ได้)

$$N(x) = Ae^{-\lambda(x-x_0)^2} \quad (13)$$

โดย x_0 คือค่าสูงสุดในช่วง (Local maximum) ในกรณี $\lambda > 0$ และให้ $A = 1$ ทำให้ค่าสูงสุดเป็น 1 เสมอสำหรับ $\forall x \forall x_0 \exists \lambda$ ซึ่งหากต้องการเทียบกับค่าสูงสุดของช่วง อาจหาได้ด้วยสมการที่ 14

$$S_N = e^{-\lambda(\frac{x-x_0}{x_0})^2} \quad (14)$$

สำหรับพารามิเตอร์ควบคุมฟังก์ชัน S_N มีเพียง λ ที่ควบคุมความโด่ง สามารถปรับได้ตามความเหมาะสมของการกระจายของข้อมูลนั้น หรืออาจใช้ $\frac{1}{2\sigma^2}$ และ $\frac{1}{\sqrt{2\pi}\sigma}$ ซึ่งเป็นค่าสัมประสิทธิ์ที่มีพารามิเตอร์ส่วนเบี่ยงเบนมาตรฐานมาช่วยได้ และ x_0 ใช้เป็นค่าสูงสุดของโค้ง (เหมาะสมที่สุด ซึ่งมีค่าเป็น 1)

7.2 ข้อได้เปรียบและข้อจำกัดของ Gaussian Function ใน Symmetric Distribution

การใช้ Normal/Non-normal Distribution มาประยุกต์หาค่าความเหมาะสม (Suitability Score) เป็นการให้ข้อมูลมีการกระจายความเหมาะสมแบบเท่ากัน และสมมาตรสองฝั่งโดยปริยาย ซึ่งปริมาณทางฟิสิกส์บางตัวแปร ไม่ได้มีความเหมาะสมและอัตราการเปลี่ยนแปลง ณ จุดข้อมูล ของความเหมาะสมเท่ากันตลอดในแนวสมมาตร อาจมีความเบ้ของความเหมาะสม ว่าข้อมูลที่อยู่ฝั่งซ้าย อาจเบ้น้อยกว่าข้อมูลที่อยู่ฝั่งขวาก็ได้ ทำให้การใช้ Symmetric Distribution อาจทำให้ค่าเพี้ยนไป ซึ่งในที่นี่ การใช้ Asymmetric Distribution อาจสื่อความหมายได้ดีกว่า และคำนวณได้แม่นยำ (ใกล้เคียงสภาพความเป็นจริง) มากกว่าแบบดังกล่าว

รวมถึงความสามารถในการประมาณค่า แบ่งเกณฑ์คิด ทดลองหา หรือทำ Curve fitting เพื่อหาความเบ้และความโด่งของโค้งจากพารามิเตอร์ที่เรากำหนด ไว้เป็นตัวแปรควบคุมได้อย่างเหมาะสมและสื่อความหมายได้ดี

7.3 การสร้างฟังก์ชัน Gaussian Suitability Score, S_N

เริ่มจากการสร้างไฟล์ของฟังก์ชันขึ้นมา โดยตั้งชื่อไฟล์ว่า `gaussianscore.m` โดยเขียนฟังก์ชันขึ้นมาพร้อมกับ argument ที่ได้กล่าวไปดังนี้

```
1 function out = gaussianscore(t,t0,sigma,lambda)
2     e = exp(1);
3     out = 1/(sqrt(2*pi*sigma))* e.^(-1/(2*sigma.^2)*lambda* ((t-t0)/t0).^2)
4 end
```

8 การนำ Beta Distribution มาประยุกต์ใช้งาน

8.1 ที่มาของฟังก์ชันหาค่าความเหมาะสม S_β

การสร้างฟังก์ชัน Beta Distribution ที่ขึ้นกับตัวแปรที่นอกเหนือจาก α และ β สามารถทำได้โดยการนำตัวแปรที่เราต้องการให้เป็น parameter (ทำเป็นตัวแปรควบคุม) และ variable (ทำเป็นตัวแปรต้น) มาจัดรูปเพื่อแทนค่าเข้าไปใน parameter เดิม ให้เป็นฟังก์ชันใหม่ ในที่นี้คือการสร้างฟังก์ชันหาค่าความเหมาะสม S_β เมื่อเปรียบเทียบค่าหรือค่าตัวแทนของชุดข้อมูลที่เก็บมาได้ เทียบกับค่าที่เหมาะสมหรือดีที่สุดค่าหนึ่ง ซึ่งข้อจำกัดของช่วง (constraint) ที่ต้องการกำหนดคือ ค่าที่ดีที่สุด (t_0) ค่าต่ำสุดของช่วงที่เร็วที่สุด (minimum, m) ค่าสูงสุดของช่วงที่เร็วที่สุด (maximum, M) และส่วนเบี่ยงเบนมาตรฐาน (standard deviation, σ) ทั้งหมดเป็น fixed parameter ส่วน variable คือค่าของข้อมูล หรือตัวแทนของชุดข้อมูลที่เก็บมาได้ (t)

กำหนดฟังก์ชันเชิงเส้น $T(t)$ ดังสมการที่ 15

$$T(t) = at + b \quad (15)$$

ข้อสังเกต 1: เมื่อเปลี่ยน b จะทำให้ Y -intercept เปลี่ยน แต่หากเปลี่ยน t เป็น ฟังก์ชันประกอบอื่น ๆ (composite function) จะทำให้ X -intercept เปลี่ยน รวมถึงการคูณค่าคงที่ไม่ใช่ค่าคง (arbitrary constant) เพื่อลด scale ของโดเมน t ลง ในที่นี้ใช้ scale เป็นพิสัยของข้อมูลที่เป็นไปได้ $\frac{1}{M-m}$ คูณเข้าไปในฟังก์ชันประกอบเป็น $T\left(\frac{t}{M-m}\right)$ จะได้ สมการที่ 16

$$T(t) = \frac{t}{M-m} + b \quad (16)$$

รวมให้ b มีค่าขยับแกน X ด้วยจะได้ สมการที่ 17

$$\begin{aligned} T(t) &= \frac{t}{M-m} + \frac{m}{m-M} \\ T(t) &= \frac{t-m}{M-m} \end{aligned} \quad (17)$$

สมการสุดท้ายที่ต้องการหาค่าความเหมาะสม S_β จะอ้างอิงจาก Beta Distribution $B(x)$ แทนค่า $T(t)$ เข้าไปเพื่อเปลี่ยน constraint จะได้ สมการที่ 18

$$S_\beta = B(T(t)) \quad (18)$$

โดย

$$B(\xi) = \frac{1}{\left[\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}\right]} \xi^{\alpha-1}(1-\xi)^{\beta-1}$$

และค่าต่ำสุดของ $B(\xi)$ มีค่าเป็น 0 จาก

$$\lim_{\xi \rightarrow m^+} B(\xi) = 0$$

และ

$$\lim_{\xi \rightarrow M^-} B(\xi) = 0$$

หากให้ $\alpha, \beta > 1$ เมื่อหาค่าสูงสุดของ $B(\xi)$ จากการหาอนุพันธ์ฟังก์ชันเทียบค่าเท่ากับศูนย์

$$\frac{\partial B(\xi)}{\partial \xi} = 0$$

แก้สมการจะได้ $B(\xi)$ มีค่าสูงสุดที่

$$\xi_0 = \frac{\alpha - 1}{\alpha + \beta - 2} \quad (19)$$

ข้อสังเกต 2: เรนจ์ของ $B(\xi)$ มีค่าสูงสุดไม่เท่ากับ 1 แต่ขึ้นกับ α และ β ทำให้จำเป็นต้องใช้วิธี Inverse scaling ผันเรนจ์จาก $[0, \xi_0]$ เป็น $[0, 1]$ ด้วยการคูณ $\frac{1}{B(\xi_0)}$ เข้าไปใน $B(\xi)$ ได้ สมการที่ 20

$$\begin{aligned} B(\xi) &= \frac{1}{B(\xi_0)} \frac{1}{\left[\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \right]} \xi^{\alpha-1} (1-\xi)^{\beta-1} \\ B(\xi) &= \frac{\left[\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \right]}{\xi_0^{\alpha-1} (1-\xi_0)^{\beta-1}} \frac{1}{\left[\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \right]} \xi^{\alpha-1} (1-\xi)^{\beta-1} \\ B(\xi) &= \frac{\xi^{\alpha-1} (1-\xi)^{\beta-1}}{\xi_0^{\alpha-1} (1-\xi_0)^{\beta-1}} \end{aligned} \quad (20)$$

เมื่อแทน $\xi_0 = \frac{\alpha-1}{\alpha+\beta-2}$ เข้าไปในสมการที่ 20 และจัดรูป จะได้ผลสุดท้ายเป็น สมการที่ 21 ดังนี้

$$\begin{aligned} B(\xi) &= \frac{(\alpha + \beta - 2)^{\alpha+\beta-2}}{(\alpha - 1)^{\alpha-1} (\beta - 1)^{\beta-1}} \xi^{\alpha-1} (1-\xi)^{\beta-1} \\ B(\xi) &= \frac{(\alpha + \beta - 2) \uparrow^2 2}{[(\alpha - 1) \uparrow^2 2] [(\beta - 1) \uparrow^2 2]} \xi^{\alpha-1} (1-\xi)^{\beta-1} \end{aligned} \quad (21)$$

เมื่อได้สมการที่ 21 แล้ว นำฟังก์ชัน $T(t)$ เข้ามาแทนค่าใน $B(\xi)$ จะได้สมการที่ 22

$$S_\beta = \frac{(\alpha + \beta - 2) \uparrow^2 2}{[(\alpha - 1) \uparrow^2 2] [(\beta - 1) \uparrow^2 2]} (T(t))^{\alpha-1} (1-T(t))^{\beta-1} \quad (22)$$

ข้อสังเกต 3: ค่าความเหมาะสม S_β ยังคงติดตัวแปร α และ β อยู่ในขณะที่ยังไม่มีตัวแปร t_0 และ σ โดยสามารถเลือกการจัดตัวแปรใดตัวแปรหนึ่งได้ แล้วผลเฉลยของ α กับ β มีความคล้ายคลึงกันมาก เนื่องจาก 2 ตัวแปรนี้แปรผกผันกันแบบสมมาตร

เราจึงสามารถแก้สมการหาว่าจุดสูงสุดของ $S_\beta = 1$ อยู่ที่ $t = t_0$ แก้สมการหาตัวแปร α ในเทอมของ t_0 ได้จาก สมการที่ 17 กับ 19 จะได้สมการที่ 23

$$\begin{aligned} \frac{\alpha - 1}{\alpha + \beta - 2} &= \frac{t_0 - m}{M - m} \\ \alpha &= \frac{(t_0 - m)(\beta - 2) + (M - m)}{M - t_0} \end{aligned} \quad (23)$$

และ β หาจากพื้นที่ใต้โค้ง Normal Distribution กับตัวแปรสมมติเพื่อปรับความโด่งของกราฟ s ให้เป็นค่าคงที่ไม่เจาะจง ซึ่งเป็นหนึ่งใน fixed parameter ของฟังก์ชันด้วย ดังสมการที่ 24

$$\beta = \frac{1}{\sqrt{2\pi}\sigma} s \quad (24)$$

8.2 โดเมนและเรนจ์ของฟังก์ชัน

8.2.1 โดเมนของตัวแปรและพารามิเตอร์คงที่

$$\alpha \in (1, +\infty)$$

$$\beta \in (1, +\infty)$$

$$t_0 \in [m, M]$$

$$s \in (\sqrt{2\pi\sigma}, +\infty)$$

8.2.2 เรนจ์ของ S_β

$$B(T(t)) \in [0, 1]$$

8.3 การสร้างฟังก์ชัน Beta Distribution Suitability Score

เริ่มจากการสร้างไฟล์ของฟังก์ชันขึ้นมา โดยตั้งชื่อไฟล์ว่า `betascore.m` โดยเขียนฟังก์ชันขึ้นมาพร้อมกับ argument ที่ได้กล่าวไปดังนี้

```
1 function out = betascore(t,t0,m,M,s,sigma)
2     e = exp(1);
3     beta = s/sqrt(2*pi*sigma);
4     alpha = ((t0-m)*(beta-2)+(M-m))/(M-t0);
5     Tex = (t-m)/(M-m);
6     o1 = (alpha+beta-2).^(alpha+beta-2)/ ...
           ((alpha-1).^(alpha-1)*(beta-1).^(beta-1));
7     o2 = Tex.^(alpha-1)*(1-Tex).^(beta-1);
8     out = o1*o2;
9 end
```

9 เนื้อหาเพิ่มเติม

9.1 สัญลักษณ์ Knuth's up-arrow Notation

สัญลักษณ์ $a \uparrow^n b$ คือ Knuth's up-arrow notation ไว้บอกการเป็น Hyperoperation, Tetration, Pentation, Hexation, ... เช่น

$$\begin{array}{lll} a \uparrow 2 = a^2 & a \uparrow^2 2 = a^a & a \uparrow^3 2 = a^{a^a} \\ a \uparrow 3 = a^3 & a \uparrow^2 3 = a^{a^a} & a \uparrow^3 3 = a^{a^{a^a}} \end{array}$$

สามารถไปศึกษาเพิ่มเติมได้ที่หัวข้อที่กล่าวไป

9.2 ตัวอย่างการสร้างฟังก์ชันคำนวณด้วย MATLAB (เก่า)

```
1 function out = calc(file)
2 for loop = 1:3
3     for i = 7:12
4         cla.figure(i)
5     end
6     dataset = csvread(file);
7     %dataset = magic(49);
8     temp = dataset(:,7);
9     humid = dataset(:,8);
10    rainfall = dataset(:,9);
11    red = dataset(:,10);
12    blue = dataset(:,11);
13    ccn = dataset(:,12);
14    carbon = dataset(:,13);
15
16    %rainfall = ANNUAL RAINFALL
17    temp1 = dataset(:,14);
18    humid1 = dataset(:,15);
19    rainfall1 = dataset(:,16);
20    red1 = dataset(:,17);
21    blue1 = dataset(:,18);
22    carbon1 = dataset(:,19);
23
24    temp2 = dataset(:,20);
25    humid2 = dataset(:,21);
26    rainfall2 = dataset(:,22);
27    red2 = dataset(:,23);
28    blue2 = dataset(:,24);
29    carbon2 = dataset(:,25);
30
31    temp3 = dataset(:,26);
32    humid3 = dataset(:,27);
33    rainfall3 = dataset(:,28);
34    red3 = dataset(:,29);
35    blue3 = dataset(:,30);
36    carbon3 = dataset(:,31);
37
38    temp4 = dataset(:,32);
39    humid4 = dataset(:,33);
40    rainfall4 = dataset(:,34);
41    red4 = dataset(:,35);
42    blue4 = dataset(:,36);
43    carbon4 = dataset(:,37);
44
45    temp5 = dataset(:,38);
46    humid5 = dataset(:,39);
47    rainfall5 = dataset(:,40);
48    red5 = dataset(:,41);
49    blue5 = dataset(:,42);
50    carbon5 = dataset(:,43);
51
52    temp6 = dataset(:,44);
53    humid6 = dataset(:,45);
54    rainfall6 = dataset(:,46);
55    red6 = dataset(:,47);
56    blue6 = dataset(:,48);
57    carbon6 = dataset(:,49);
58
59    r(1,1) = corr(temp,temp1);
60    r(1,2) = corr(temp,temp2);
61    r(1,3) = corr(temp,temp3);
62    r(1,4) = corr(temp,temp4);
63    r(1,5) = corr(temp,temp5);
64    r(1,6) = corr(temp,temp6);
65
66    r(2,1) = corr(humid,humid1);
67    r(2,2) = corr(humid,humid2);
68    r(2,3) = corr(humid,humid3);
69    r(2,4) = corr(humid,humid4);
70    r(2,5) = corr(humid,humid5);
71    r(2,6) = corr(humid,humid6);
72
73    r(3,1) = corr(rainfall,rainfall1);
74    r(3,2) = corr(rainfall,rainfall2);
75    r(3,3) = corr(rainfall,rainfall3);
76    r(3,4) = corr(rainfall,rainfall4);
77    r(3,5) = corr(rainfall,rainfall5);
78    r(3,6) = corr(rainfall,rainfall6);
79
80    r(4,1) = corr(red,red1);
81    r(4,2) = corr(red,red2);
82    r(4,3) = corr(red,red3);
```

```

83 r(4,4) = corr(red,red4);
84 r(4,5) = corr(red,red5);
85 r(4,6) = corr(red,red6);
86
87 r(5,1) = corr(blue,blue1);
88 r(5,2) = corr(blue,blue2);
89 r(5,3) = corr(blue,blue3);
90 r(5,4) = corr(blue,blue4);
91 r(5,5) = corr(blue,blue5);
92 r(5,6) = corr(blue,blue6);
93
94 r(6,1) = corr(carbon,carbon1);
95 r(6,2) = corr(carbon,carbon2);
96 r(6,3) = corr(carbon,carbon3);
97 r(6,4) = corr(carbon,carbon4);
98 r(6,5) = corr(carbon,carbon5);
99 r(6,6) = corr(carbon,carbon6);
100
101 %temp
102
103 dat = dataset(:,7);
104 dat1 = dataset(:,14);
105 dat2 = dataset(:,20);
106 dat3 = dataset(:,26);
107 dat4 = dataset(:,32);
108 dat5 = dataset(:,38);
109 dat6 = dataset(:,44);
110
111 p = polyfit(dat1,dat,1);
112 f = polyval(p,dat1);
113
114 hold on
115
116 figure(7)
117 plot(dat1,f,'--r')
118 plot(dat,dat1)
119 scatter(dat,dat1)
120 plot(dat2,f,'--r')
121 plot(dat,dat2)
122 scatter(dat,dat2)
123 plot(dat3,f,'--r')
124 plot(dat,dat3)
125 scatter(dat,dat3)
126 plot(dat4,f,'--r')
127 plot(dat,dat4)
128 scatter(dat,dat4)
129 plot(dat5,f,'--r')
130 plot(dat,dat5)
131 scatter(dat,dat5)
132 plot(dat6,f,'--r')
133 plot(dat,dat6)
134 scatter(dat,dat6)
135 plot(dat,dat)
136
137
138 hold off
139
140 %humid
141
142 dat = dataset(:,8);
143 dat1 = dataset(:,15);
144 dat2 = dataset(:,21);
145 dat3 = dataset(:,27);
146 dat4 = dataset(:,33);
147 dat5 = dataset(:,39);
148 dat6 = dataset(:,45);
149
150 p = polyfit(dat1,dat,1);
151 f = polyval(p,dat1);
152
153 hold on
154
155 figure(8)
156 plot(dat1,f,'--r')
157 plot(dat,dat1)
158 scatter(dat,dat1)
159 plot(dat2,f,'--r')
160 plot(dat,dat2)
161 scatter(dat,dat2)
162 plot(dat3,f,'--r')
163 plot(dat,dat3)
164 scatter(dat,dat3)
165 plot(dat4,f,'--r')
166 plot(dat,dat4)
167 scatter(dat,dat4)
168 plot(dat5,f,'--r')
169 plot(dat,dat5)
170 scatter(dat,dat5)
171 plot(dat6,f,'--r')
172 plot(dat,dat6)
173 scatter(dat,dat6)
174 plot(dat,dat)
175
176 hold off
177
178 %rainfall
179
180 dat = dataset(:,9);
181 dat1 = dataset(:,16);
182 dat2 = dataset(:,22);
183 dat3 = dataset(:,28);
184 dat4 = dataset(:,34);
185 dat5 = dataset(:,40);
186 dat6 = dataset(:,46);
187
188 p = polyfit(dat1,dat,1);
189 f = polyval(p,dat1);
190
191 hold on
192
193 figure(9)
194 plot(dat1,f,'--r')
195 plot(dat,dat1)
196 scatter(dat,dat1)
197 plot(dat2,f,'--r')
198 plot(dat,dat2)

```

```

199 scatter(dat,dat2)
200 plot(dat3,f,'--r')
201 plot(dat,dat3)
202 scatter(dat,dat3)
203 plot(dat4,f,'--r')
204 plot(dat,dat4)
205 scatter(dat,dat4)
206 plot(dat5,f,'--r')
207 plot(dat,dat5)
208 scatter(dat,dat5)
209 plot(dat6,f,'--r')
210 plot(dat,dat6)
211 scatter(dat,dat6)
212 plot(dat,dat)
213
214 hold off
215
216 %red
217 dat = dataset(:,10);
218 dat1 = dataset(:,17);
219 dat2 = dataset(:,23);
220 dat3 = dataset(:,29);
221 dat4 = dataset(:,35);
222 dat5 = dataset(:,41);
223 dat6 = dataset(:,47);
224
225 p = polyfit(dat1,dat,1);
226 f = polyval(p,dat1);
227
228 hold on
229
230 figure(10)
231 plot(dat1,f,'--r')
232 plot(dat,dat1)
233 scatter(dat,dat1)
234 plot(dat2,f,'--r')
235 plot(dat,dat2)
236 scatter(dat,dat2)
237 plot(dat3,f,'--r')
238 plot(dat,dat3)
239 scatter(dat,dat3)
240 plot(dat4,f,'--r')
241 plot(dat,dat4)
242 scatter(dat,dat4)
243 plot(dat5,f,'--r')
244 plot(dat,dat5)
245 scatter(dat,dat5)
246 plot(dat6,f,'--r')
247 plot(dat,dat6)
248 scatter(dat,dat6)
249 plot(dat,dat)
250
251 hold off
252
253 %blue
254 dat = dataset(:,11);
255 dat1 = dataset(:,18);
256 dat2 = dataset(:,24);

```

```

257 dat3 = dataset(:,30);
258 dat4 = dataset(:,36);
259 dat5 = dataset(:,42);
260 dat6 = dataset(:,48);
261
262 p = polyfit(dat1,dat,1);
263 f = polyval(p,dat1);
264
265 hold on
266
267 figure(11)
268 plot(dat1,f,'--r')
269 plot(dat,dat1)
270 scatter(dat,dat1)
271 plot(dat2,f,'--r')
272 plot(dat,dat2)
273 scatter(dat,dat2)
274 plot(dat3,f,'--r')
275 plot(dat,dat3)
276 scatter(dat,dat3)
277 plot(dat4,f,'--r')
278 plot(dat,dat4)
279 scatter(dat,dat4)
280 plot(dat5,f,'--r')
281 plot(dat,dat5)
282 scatter(dat,dat5)
283 plot(dat6,f,'--r')
284 plot(dat,dat6)
285 scatter(dat,dat6)
286 plot(dat,dat)
287
288 hold off
289
290 %carbon
291 dat = dataset(:,13);
292 dat1 = dataset(:,19);
293 dat2 = dataset(:,25);
294 dat3 = dataset(:,31);
295 dat4 = dataset(:,37);
296 dat5 = dataset(:,43);
297 dat6 = dataset(:,49);
298
299 p = polyfit(dat1,dat,1);
300 f = polyval(p,dat1);
301
302 hold on
303
304 figure(12)
305 plot(dat1,f,'--r')
306 plot(dat,dat1)
307 scatter(dat,dat1)
308 plot(dat2,f,'--r')
309 plot(dat,dat2)
310 scatter(dat,dat2)
311 plot(dat3,f,'--r')
312 plot(dat,dat3)
313 scatter(dat,dat3)
314 plot(dat4,f,'--r')

```



```

315 plot(dat,dat4)
316 scatter(dat,dat4)
317 plot(dat5,f,'--r')
318 plot(dat,dat5)
319 scatter(dat,dat5)
320 plot(dat6,f,'--r')
321 plot(dat,dat6)
322 scatter(dat,dat6)
323 plot(dat,dat)
324
325 hold off
326
327 for i = 1:6
328     r(7,i) = (r(1,i)+r(2,i)+r(3,i)+ ...
329             r(4,i)+r(5,i)+r(6,i))/6;
330
331 for i = 1:7
332     for j = 1:6
333         r(i,j) = r(i,j).^2;
334     end
335 end
336
337 r = 100*r;
338
339 csvwrite('G:\Project CANSAT18 ...
340          Final Round\MATLAB\export.csv',r)
341
342 %type('F:\Project CANSAT18 Final ...
343       Round\MATLAB\export.csv')
344 end
345
346 end

```

บรรณานุกรม

- [1] David Houcque. (2005). *Introduction to **MATLAB** for Engineering Students*. Northwestern University.
- [2] Johnathan Mun. (2008). *Advanced Analytical Models, Understanding and Choosing the Right Probability Distributions*. New York