

Linear Regression

Petchara Pattarakijwanich

Introduction to Data Science, 9 September 2022

Goal of this week

- Simple Linear Regression
 - Least Squared Fitting
 - Significance and p-value
 - Multiple Regression
 - Interaction and Non-linear Terms
 - Some technical details
-
- (Logistic Regression)

Regression

- $y = \beta_0 + \beta_1 x$ \Rightarrow Simple
- $y = \beta_0 + \beta_1 x + \beta_2 x^2$ \Rightarrow Quadratic
- $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$ \Rightarrow Cubic

(Ridge, Lasso, subset selection)

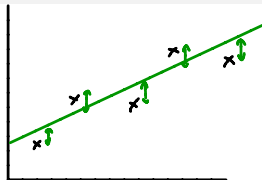
Simple Linear Regression

$$y = \beta_0 + \beta_1 x + \epsilon$$

← noise

ค่า y เมื่อ $x = 0$

x เพิ่มขึ้น 1 หน่วย y เพิ่มขึ้นเท่าไร



หา β_0, β_1 ที่ "ดีที่สุด" \Rightarrow ทำให้ $RSS = \sum_i (y - \beta_0 - \beta_1 x)^2$

น้อยที่สุด

Residual sum of Square

$$\chi^2 = \sum_i \frac{(y_{data} - y_{model})^2}{\sigma^2}$$

↘
อธิบาย

$$RSS = \sum_i (y_{data} - y_{model})^2$$

$$RSS = \chi^2 \text{ ถ้า } \sigma \text{ คง}$$

Linear Least Squared Fitting

Let β_0, β_1 minimize $\sum_i (y_i - \beta_0 - \beta_1 x_i)^2$ (minimize)

$$RSS = \sum_i (y_i - \beta_0 - \beta_1 x_i)^2$$

$$= (y_1 - \beta_0 - \beta_1 x_1)^2 + (y_2 - \beta_0 - \beta_1 x_2)^2 + \dots + (y_n - \beta_0 - \beta_1 x_n)^2$$

\uparrow x_1, \dots, x_n & y_1, \dots, y_n are given = data

\Rightarrow RSS is a function of β_0, β_1

$$RSS = a\beta_0^2 + b\beta_1^2 + c\beta_0\beta_1 + d\beta_0 + e\beta_1 + f$$

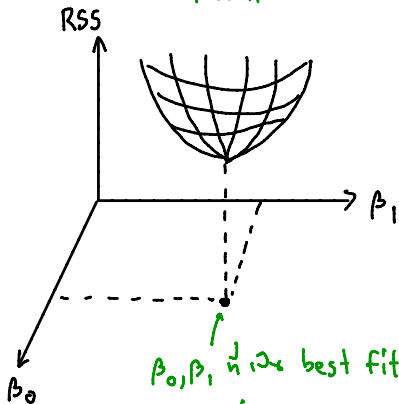
\uparrow \uparrow
 n $\sum (x_i)^2$
 $\nwarrow \nearrow$
 independent

Linear Least Squared Fitting

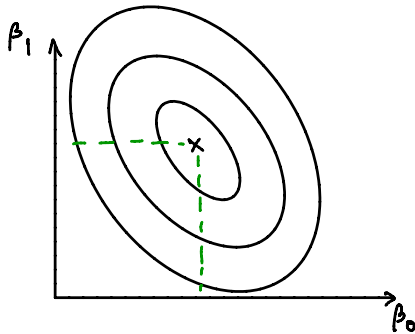
$$RSS = a\beta_0^2 + b\beta_1^2 + c\beta_0\beta_1 + d\beta_0 + e\beta_1 + f$$

n
 $\sum(x_i)^2$
ข้อมูล

Paraboloid 1/2 3 มิติ



β_0, β_1 ให้ best fit
RSS หายไป



Linear Least Squared Fitting

$$RSS = \sum_i (y_i - \beta_0 - \beta_1 x_i)^2$$

หาค่าของ $\frac{\partial(RSS)}{\partial \beta_0} = 0$, $\frac{\partial(RSS)}{\partial \beta_1} = 0$

$$\frac{\partial RSS}{\partial \beta_0} = \sum_i 2(y_i - \beta_0 - \beta_1 x_i)(-1) = -2 \left[\sum_i y_i - \beta_0 \overset{=n}{\sum_i 1} - \beta_1 \sum_i x_i \right] = 0$$

$$\frac{\partial RSS}{\partial \beta_1} = \sum_i 2(y_i - \beta_0 - \beta_1 x_i)(-x_i) = -2 \left[\sum_i x_i y_i - \beta_0 \sum_i x_i - \beta_1 \sum_i x_i^2 \right] = 0$$

$$\beta_0 n + \beta_1 \sum_i x_i = \sum_i y_i$$

$$\beta_0 \sum_i x_i + \beta_1 \sum_i x_i^2 = \sum_i x_i y_i$$

$$\Rightarrow \begin{pmatrix} n & \sum x \\ \sum x & \sum x^2 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} = \begin{pmatrix} \sum y \\ \sum xy \end{pmatrix}$$

$$\Rightarrow \text{แก้สมการ 2 สมการ 2 ตัว} (\beta_0, \beta_1)$$

Linear Least Squared Fitting

aimou

$$\beta_1 = \frac{n \sum_i x_i y_i - (\sum_i x_i)(\sum_i y_i)}{n \sum_i x_i^2 - (\sum_i x_i)^2}$$

$$\beta_0 = \frac{1}{n} \left[\sum_i y_i - \beta_1 \sum_i x_i \right]$$

Significance and p-value

$$y = \beta_0 + \beta_1 x + \epsilon$$

y ខឹង x ឬ ទេ? $\iff \beta_1 \neq 0$ ឬ ទេ?

Hypothesis test

- Null Hypothesis $\beta_1 = 0$
- Alternative hypothesis $\beta_1 \neq 0$



ឬ Probability ក្នុងការសង្កេតឃើញ $\beta_1 = 0$ [F-statistics]



ឬ probability ក្នុងការសង្កេតឃើញ $\beta_1 \neq 0$ [p-value]
ដែល ជាធម្មតា ត្រូវបានប្រើប្រាស់

Multiple Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon \quad \leftarrow \text{residual } n+1 \text{ value}$$

β_0 คือค่า y เมื่อ $x_1 = x_2 = \dots = x_n = 0$

β_1 คือ Δy เมื่อ $\Delta x_1 = 1$ เมื่อ x ที่เหลือคงที่

หา $\beta_0, \beta_1, \dots, \beta_n$ ที่ทำให้ $\text{RSS} = \sum (y - \beta_0 - \beta_1 x_1 - \dots - \beta_n x_n)^2$ มีค่าน้อยสุด

$$\Rightarrow \frac{\partial \text{RSS}}{\partial \beta_i} = 0 \Rightarrow n+1 \text{ สมการ } n+1 \text{ ตัวแปร}$$

\Rightarrow แก้หา β ได้ $[\beta_0, \beta_1, \dots, \beta_n]$

Interaction Term and Non-linear Term

• Interaction term: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2$

$$y = \beta_0 + (\beta_1 + \beta_{12} x_2) x_1 + \beta_2 x_2 \quad [\text{သတ်မှတ်သော } x_1 \text{ ၏ပေါ်တွင် } x_2]$$

$$y = \beta_0 + \beta_1 x_1 + (\beta_2 + \beta_{12} x_1) x_2 \quad [\text{သတ်မှတ်သော } x_2 \text{ ၏ပေါ်တွင် } x_1]$$

• Non-linear term: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \alpha_1 x_1^2 + \alpha_2 x_2^2$

မူလ linear ၏ x ၏ပေါ်တွင် linear ၏ β

\Rightarrow ဤ β ၏ပေါ်တွင် RSS ၏ပေါ်တွင် \Rightarrow သတ်မှတ်သော သတ်မှတ်သော β

Technical Problems

- Parameter normalization?

Ex. x_1 သတ်၍ $\approx 10^{-5}$
 x_2 သတ်၍ $\approx 10^{20}$ } \Rightarrow overflow
underflow
Round-off error } ကိန်း normalization

- Error bars?

scale x_i လွှဲ၍ $\bar{x}_i = 0$, $\text{var}(x_i) = 1$

$$\text{RSS} = \sum (y - y_{\text{model}})^2 \Rightarrow \chi^2 = \sum \frac{(y - y_{\text{model}})^2}{\sigma^2}$$

- Qualitative input?

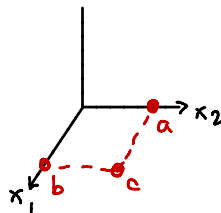
2 ကိန်း $\{a, b\}$

$$x = \begin{cases} 0, & a \\ 1, & b \end{cases}$$

3 ကိန်း $\{a, b, c\}$

$$x_1 = \begin{cases} 0, & a \\ 1, & \text{not } a \end{cases}$$

$$x_2 = \begin{cases} 0, & b \\ 1, & \text{not } b \end{cases}$$



Technical Problems

- Non-linearity

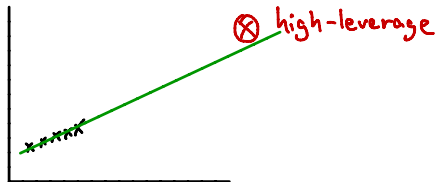
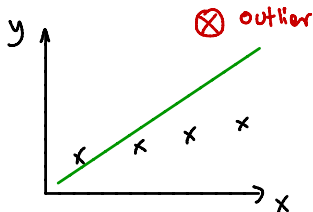
- many Non-linear and interaction terms
 $\{x_i^2, x_i^3, \dots\}$ $\{x_i x_j, \dots\}$

- Collinearity

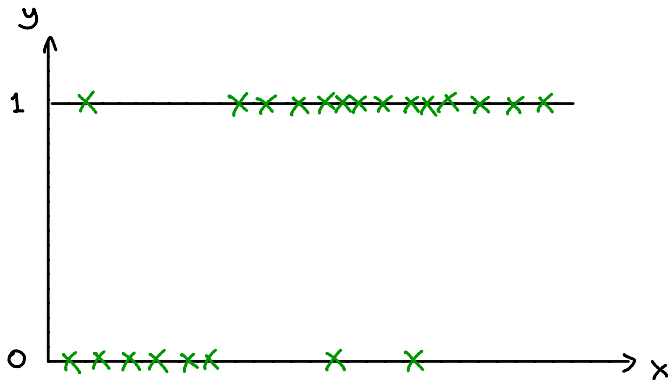
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

ถ้า x_1, x_2 มีความสัมพันธ์กัน? $x_1 = kx_2 \Rightarrow$ ใช้เวลา information
เหมือนกันว่า x_1, x_2 ใด ๆ

- Outliers & High-leverage points



Logistic Regression Example

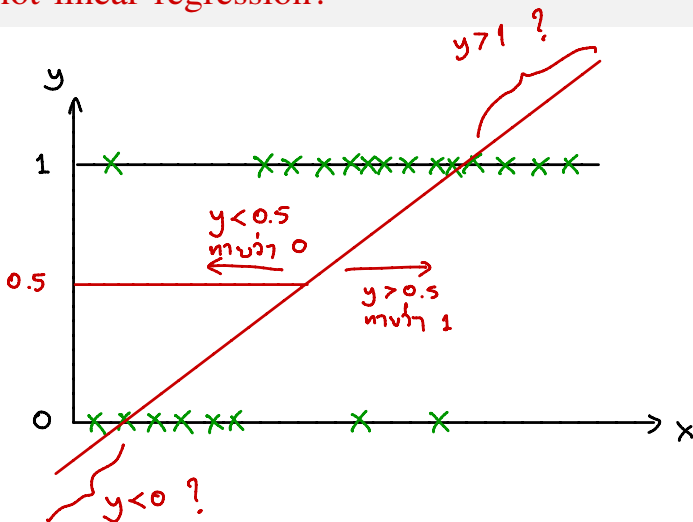


x = 1 หรือ 0

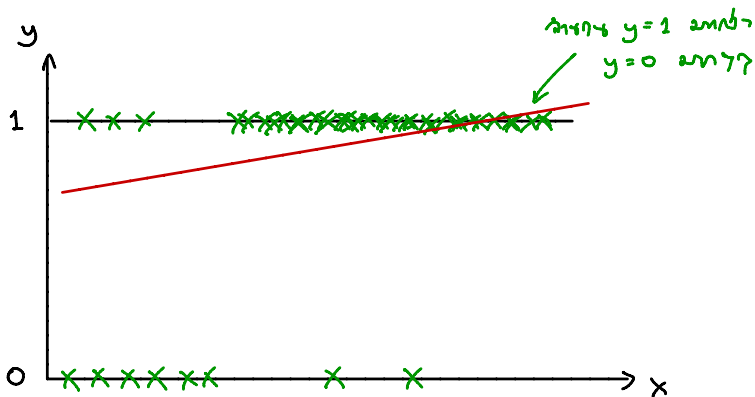
y = 1 หรือ 0 = $\begin{cases} 1, \text{ใช่} \\ 0, \text{ไม่ใช่} \end{cases}$

ตัวอย่าง ทำตาม
ถ้า y กับ x ?

Why not linear regression?

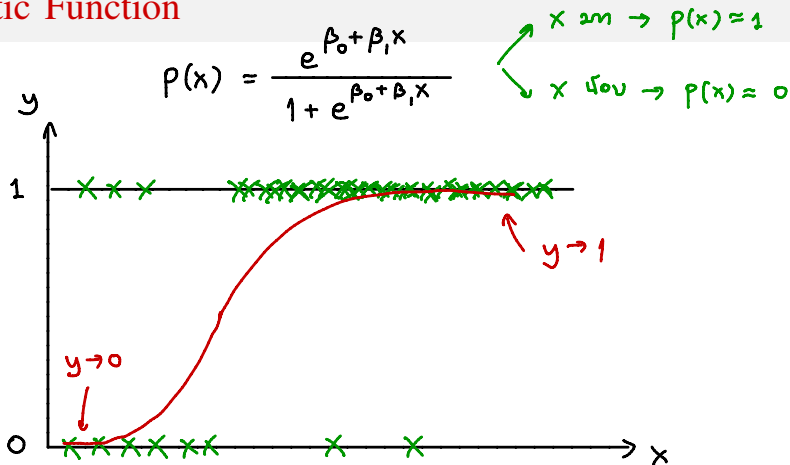


Why not linear regression?



Linear regression fails

Logistic Function

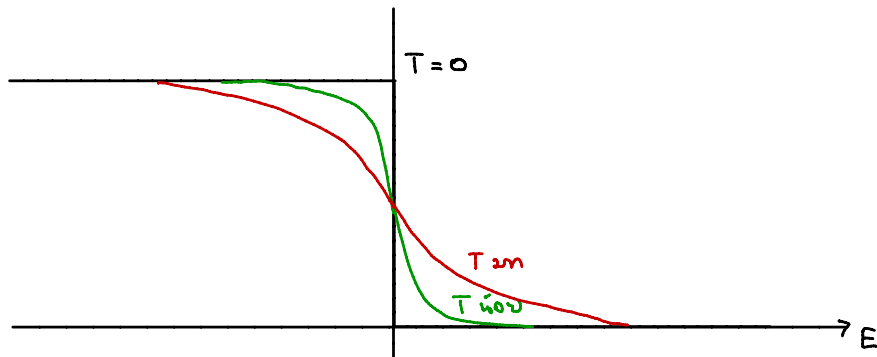


คุณสมบัติ \rightarrow $p(x) \in [0, 1]$
monotone

Logistic Function

Fermi-Dirac Distribution

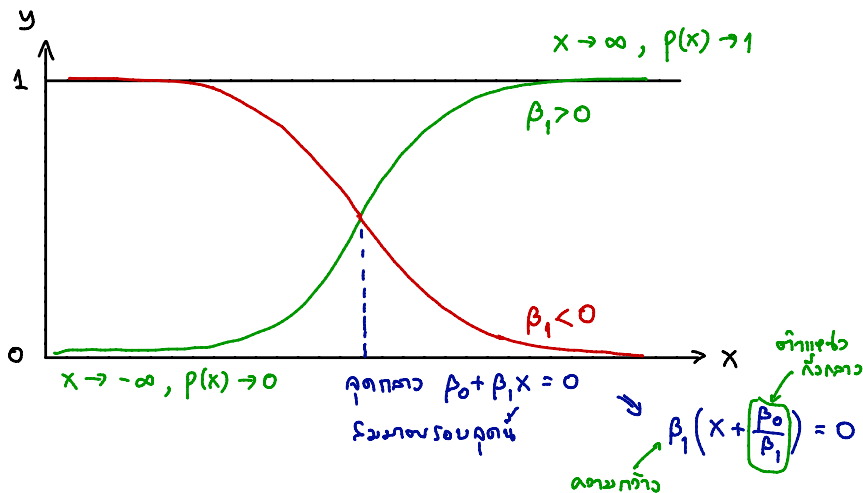
$$f(E) = \frac{e^{-(E-E_F)/kT}}{1 + e^{-(E-E_F)/kT}}$$



$$f(E) = \begin{cases} \rightarrow 0 & \text{เมื่อ } E \gg E_F \quad (E \rightarrow \infty) \\ \rightarrow 1 & \text{เมื่อ } E \ll E_F \quad (E \rightarrow -\infty) \\ 0.5 & \text{เมื่อ } E = E_F \end{cases}$$

Logistic Function

$$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$



Odds and logit function

$$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} = \frac{e^{\alpha}}{1 + e^{\alpha}}$$

$$[\ln p \propto \ln \frac{p}{1-p}]$$

$$(1 + e^{\alpha}) p = e^{\alpha}$$

$$p = (1-p) e^{\alpha}$$

$$e^{\alpha} = \frac{p}{1-p}$$

$$\alpha = \beta_0 + \beta_1 x = \ln \left(\frac{p}{1-p} \right)$$

logit(p) is
linear function of x

logit function

logistic unit

ကျွန်ုပ်တို့ "odd"

တစ်ခုခုဖြစ်ရန် အားသာမှုရှိခြင်း
↙ event ခုဖြစ်ရန် ချီခြင်း

$$\text{odd} = \frac{p}{1-p}$$

⇒ အားသာမှု logistic function
ရှိခြင်း probability နှင့် $y=1$

Logistic Regression

ພົບ β_0, β_1 ດ້ວຍ "Maximum Likelihood"

likelihood \mathcal{L} = Probability ທີ່ຈະໄດ້ຜົນຖອດຢ່າງນີ້ ຖ້າຂໍ້ມູນທີ່ໃຫ້ β_0, β_1
= $P(\text{data} \mid \text{model})$

$$= P(\text{ອາດເປັນ } 1) \times P(\text{ອາດເປັນ } 2) \times \dots \times P(\text{ອາດເປັນ } n)$$

ຂໍ້ມູນ $y_1 = 1$ $y_2 = 0$... $y_n = 0$

\uparrow \uparrow

$P(\text{ອາດເປັນ } 1) = P(x_1)$ $P(\text{ອາດເປັນ } 2) = 1 - P(x_2)$

$$\mathcal{L} = \prod_{\substack{\text{ທັງ } i \\ y_i = 1}} P(x_i) \prod_{\substack{\text{ທັງ } i \\ y_i = 0}} (1 - P(x_i))$$

\mathcal{L} ເປັນຜົນຄູນຂອງ β_0, β_1
 \Rightarrow ພົບ β_0, β_1 ທີ່ \mathcal{L} ມີຄ່າສູງ

Logistic Regression

$$\mathcal{L} = \prod_{\substack{n \text{ i } n \\ y_i = 1}} p(x_i) \prod_{\substack{n \text{ i } n \\ y_i = 0}} (1 - p(x_i))$$

likelihood = 1

for $i = 1, \dots, n_{\text{data}}$
if $y_i == 0$
likelihood $\times = p(x_i)$
if $y_i == 1$
likelihood $\times = 1 - p(x_i)$

$$\mathcal{L} = \prod_{n \text{ i } n} p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$$

→ likelihood = 1
for $i = 1, \dots, n_{\text{data}}$
likelihood $\times = p^{y_i} (1-p)^{1-y_i}$

หาค่าของ $\ln \mathcal{L}$ แทน \mathcal{L}

\mathcal{L} มาก $\Rightarrow \ln \mathcal{L}$ มาก