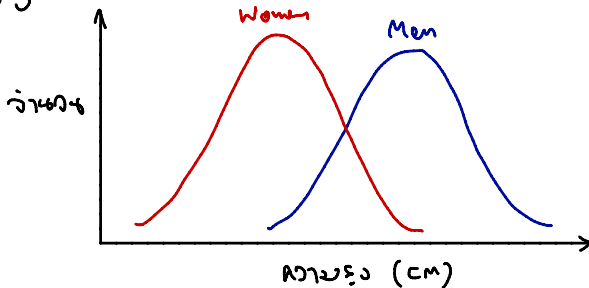


Model Selection

Petchara Pattarakijwanich

Introduction to Data Science, 21 October 2022

HW 3



M : $x_{m1}, x_{m2}, \dots, x_{m500}$

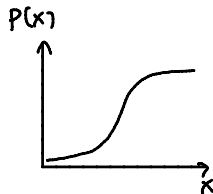
W : $x_{w1}, x_{w2}, \dots, x_{w500}$

ถ้า β_0, β_1
ห้ของมาบโตะจะได้ผลลัพธ์?

Logistic function :

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

$$p(W \text{ ห้ความสูง } x) = 1 - p(x)$$



$$\mathcal{L} = \text{probability of } \hat{y} \text{ given } x \text{ and } \beta_0, \beta_1 \text{ data}$$

$$= \underbrace{\left[p(x_{m1}) p(x_{m2}) p(x_{m3}) \dots p(x_{m508}) \right]}_{\text{prob of } \hat{y} \text{ Men}} \cdot$$

$$\underbrace{\left[(1 - p(x_{w1})) (1 - p(x_{w2})) \dots (1 - p(x_{w500})) \right]}_{\text{prob of } \hat{y} \text{ Women}}$$

$$= \prod_{\text{men}} p(x_{mi}) \prod_{\text{women}} (1 - p(x_{wi}))$$

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

$$\mathcal{L} = \mathcal{L}(\beta_0, \beta_1)$$

\Rightarrow Optimize wrt β_0, β_1 to \mathcal{L} max

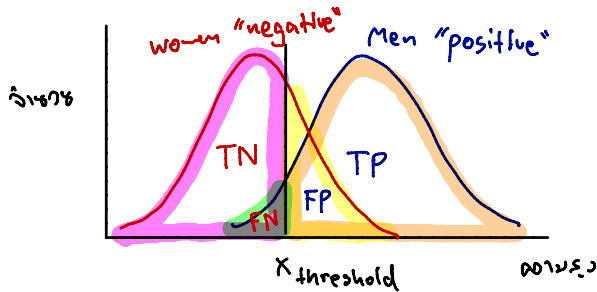
$$\mathcal{L} = \prod_{\text{men}} p(x_{mi}) \prod_{\text{women}} (1 - p(x_{wi})) \quad \leftarrow \begin{array}{l} \text{underflow} \\ \text{ถ้าคูณกันแล้ว} \end{array}$$

$$\ln \mathcal{L} = \sum_{\text{men}} \ln p(x_{mi}) + \sum_{\text{women}} \ln (1 - p(x_{wi})) \quad \leftarrow \begin{array}{l} \text{saved} \\ \text{underflow} \end{array}$$

$$-\ln \mathcal{L} = -\sum_{\text{men}} \ln p(x_{mi}) - \sum_{\text{women}} \ln (1 - p(x_{wi}))$$

↑
 $-\ln \mathcal{L}$ มีค่าลดลงเมื่อ \mathcal{L} มีค่ามากขึ้น

$$\Rightarrow \text{we } \beta_0, \beta_1 \text{ will minimize } -\ln \mathcal{L}$$



man Men if $x \geq x_{\text{threshold}}$

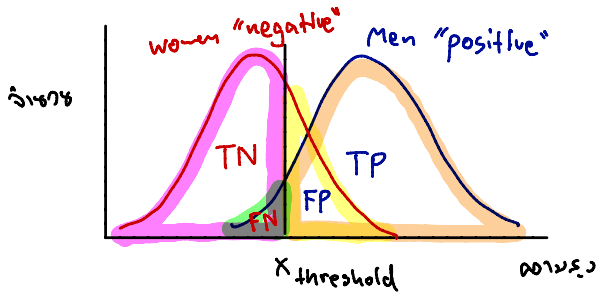
women if $x < x_{\text{threshold}}$

TP = density of men x_+

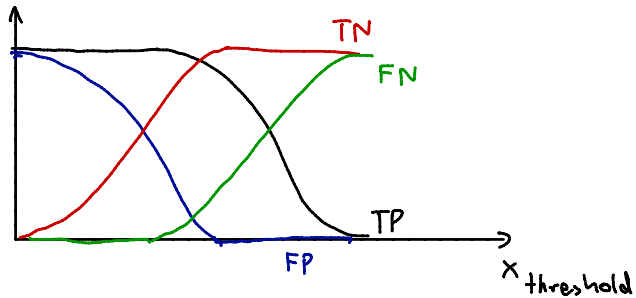
TN = density of women x_+

FP = density of men x_+

FN = density of women x_+



TP, TN
FP, FN



Goal of this week

- Standard Linear Regression and Its Limits
- Model Selection Methods
 - Mallow's C_p
 - Akaike Information Criteria (AIC)
 - Bayesian Information Criteria (BIC)
 - (Reduced χ^2)
- Subset Selection
 - Best Subset Selection
 - Forward/Backward Step-wise Subset Selection

Standard Linear Regression and Its Limits

input : $\{x_1, x_2, \dots, x_p\}$ output : y

Linear Regression :

$$y_{\text{model}} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p \\ (+ \text{ interaction terms} + \text{ non-linear terms})$$

$$\Rightarrow \text{we } \beta_0, \beta_1, \dots, \beta_p \text{ minimize } \text{RSS} = \sum (y - y_{\text{model}})^2 \text{ over } \beta \\ \left[\text{over } \text{RSS}, \frac{\partial \text{RSS}}{\partial \beta_0} = 0, \frac{\partial \text{RSS}}{\partial \beta_1} = 0, \dots \Rightarrow \text{we find } \rightarrow \text{minimize} \right]$$

\Rightarrow we statistical significance of β_i too

Standard Linear Regression and Its Limits

สมมติว่าเรามี y หนึ่งตัว x_1 หนึ่งตัว (ไม่มี x_i)

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

\Rightarrow หนึ่งตัว $\beta_2, \beta_3, \dots, \beta_p \approx 0$
และไม่มี statistical significance

} Sum of squares
ลดน้อยลง

ปัญหา :

$n \approx p$ หนึ่งตัว แล้ว variance หมด

$n < p$ หนึ่งตัว

การคำนวณ $n \gg p$
(ตัวอย่างเช่น 100-50)

Standard Linear Regression and Its Limits

precision medicine
(personalized)

Genome sequence \Rightarrow 1 Gene $\approx 500,000$ bits
(binary $\approx 500,000$ bits) **input**

Genotype : ≈ 1000 bits **output**

Genotype : ≈ 1000 bits

$$n \approx 1000, \quad p \approx 500,000$$

Standard Linear Regression and Its Limits

วิธีแก้ Linear Regression

① แก้ Linear Regression

- Subset Selection
- Dimensionality Reduction
- Regularization (Ridge/Lasso)

② แก้ (non-linear)

- Tree-based
- Support Vector Machine
- Neural Network

Model Selection

มีตัวอย่างข้อมูล 3 โมเดล

$$y_1 = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

$$y_2 = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2$$

$$y_3 = \beta_0 + \beta_1 x_1$$

Model Selection

ถ้าหากมีข้อมูล / จำนวนตัวแปร ที่ต่างกัน

⇒ โมเดลไหน "ดีที่สุด" หน่อย?

MSE_{training} ของโมเดล
โมเดลที่ขึ้นชื่อ overfit

- ทดสอบกับ test data ⇒ MSE_{test} ที่น้อยคือ "ดีที่สุด" หน่อย
- Cross Validation (ใช้ข้อมูล test data 1 ชุด)
- พยายาม Estimate MSE_{test} จาก MSE_{training} ของหลายๆ โมเดล

Model Selection

เราต้อง Estimate MSE_{test} on $MSE_{training}$ ดูหุญจ้ มาฝาก

$$MSE_{test} = \underbrace{MSE_{training}} + \text{correction}$$

$MSE_{test} > MSE_{training}$ หมายความว่า overfit

วิธีแก้ไข Correction

- Likelihood {
 - Mallows's C_p
 - Akaike Information Criteria (AIC)
- Posterior {
 - Bayesian Information Criteria (BIC)
- Ad hoc {
 - Reduced χ^2

Mallow's C_p

Residual Sum of Squares

$$= \sum (y - y_{\text{model}})^2$$

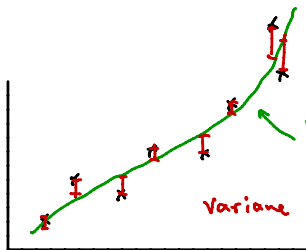
(for training data)

variance of residual
(အကွက်ချွန်ချွန်ချွန်ချွန်)
or noise variance

$$C_p = \frac{1}{n} (RSS + 2d\sigma^2)$$

ချွန်ချွန်ချွန်ချွန်

ချွန်ချွန် free parameter
ချွန်ချွန်



variance of residual = σ^2

Mallow's C_p

$$C_p = \frac{1}{n} (RSS + 2d\sigma^2)$$
$$= \frac{RSS}{n} + \frac{2d\sigma^2}{n}$$

$$C_p = \text{MSE}_{\text{training}} + \underbrace{\frac{2d\sigma^2}{n}}_{\text{correction}}$$

(subtract free parameter)

C_p ใกล้เคียง \Rightarrow โมเดลดี "พอๆ"

C_p ใกล้เคียง \nearrow โมเดลดี (MSE 600)
 \searrow free parameter น้อยเกินไป (d 600)

Mallow's C_p

$$\text{ឆ្លងកាត់រវាង } C_p \Rightarrow C'_p = \frac{RSS}{\sigma^2} + 2d - n$$

\Rightarrow តើ C_p ល្អបំផុត C'_p ល្អបំផុតដែរឬទេ

Akaike Information Criteria (AIC)

$$AIC = \frac{1}{n\sigma^2} (RSS + 2d\sigma^2)$$

$$= \frac{1}{\sigma^2} \left[\frac{RSS}{n} + \underbrace{\frac{2d\sigma^2}{n}} \right]$$

$$AIC = \frac{1}{\sigma^2} C_p$$

→ n ၁၀၀၀ free parameter

$n \gg p$ ၇၂၆၇၁၀၀

$n \approx p$ ၆၇၁၀၀

AIC ခံနိုင်ရည် \Rightarrow အကဲဖြတ် "မိုဒယ်"

Bayesian Information Criteria (BIC)

$$\text{BIC} = \frac{1}{n\sigma^2} \left(\text{RSS} + d\sigma^2 \ln n \right)$$

$$\text{BIC} = \frac{1}{\sigma^2} \left(\frac{\text{RSS}}{n} + d\sigma^2 \frac{\ln n}{n} \right)$$

$$= \frac{1}{\sigma^2} \left(\text{MSE} + d\sigma^2 \frac{\ln n}{n} \right)$$

$$\left[\text{AIC} = \frac{1}{\sigma^2} \left(\text{MSE} + \frac{2d\sigma^2}{n} \right) \right]$$

$\text{BIC} > \text{AIC}$ ถ้า $n \gtrsim 7-8$

\Rightarrow BIC จะเลือกโมเดลที่ AIC

\Rightarrow ถ้าใช้ BIC จะเลือกโมเดลที่ซับซ้อนน้อยกว่า

Reduced χ^2

$$\chi^2 = \sum \left(\frac{y - y_{\text{model}}}{\Delta y} \right)^2 \approx \text{RSS}$$

$$\text{Reduced } \chi^2 = \frac{\chi^2}{\text{dof}} = \frac{\chi^2}{n - d}$$