

Checking Predictive Power

2/24/22

Response: $Y = \begin{cases} 1 \\ 0 \end{cases}$ Prediction: $\hat{Y} = \begin{cases} 1 \\ 0 \end{cases}$

Model gives us $\hat{\pi} \in [0, 1]$.

Cut-off value: $\hat{Y} = \begin{cases} 1 & \hat{\pi} > \pi_0 \\ 0 & \hat{\pi} \leq \pi_0 \end{cases}$

For some chosen π_0 . Common choices:

- $\pi_0 = 0.5$

- $\pi_0 =$ sample proportion of $Y=1$

Classification Table:

		Predicted	
		$\hat{Y}=1$	$\hat{Y}=0$
Observed	$Y=1$	a	b
	$Y=0$	c	d

(Note: In the original image, 'a' is circled in blue, 'b' has a red double arrow pointing to it, 'c' has a red double arrow pointing to it, and 'd' is circled in blue. There are also blue double arrows pointing to 'c' and 'd'.)

Sensitivity: $P(\hat{Y}=1 | Y=1)$

Estimate: $\frac{a}{a+b}$

Specificity: $P(\hat{Y}=0 | Y=0)$

Estimate: $\frac{d}{c+d}$

Better - split data

① Training \rightarrow fit model

② Testing \rightarrow classification above

Warning:

- fit the model on the same data we're trying to predict

- the same data we're trying to predict

predict

\rightarrow not fair assessment of predictive power.

Limitations : ① Sensitive to our choice of π_0 .

② Sensitive to the relative times that $Y=1$ & $Y=0$ in the data

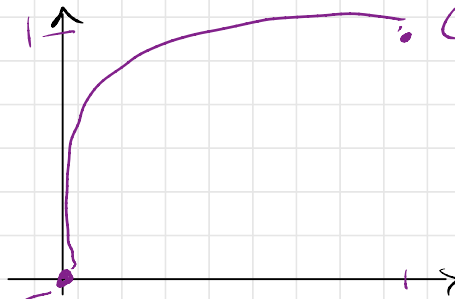
③ Loss of information classifying a probability $\hat{\pi}$ to 0, 1.

ROC - Receiver operating curves (characteristic)

- Calculate sensitivity & specificity for all $\pi_0 \in [0, 1]$

$$P(\hat{Y}=1 | Y=1)$$

Sensitivity
= true
positive rate
(want high)



$$\hat{Y} = \begin{cases} 1 & \hat{\pi} > 0 \\ 0 & \hat{\pi} \leq 0 \end{cases} \Rightarrow \hat{Y} = 1$$

$$\Rightarrow \hat{Y} = \begin{cases} 1 & \hat{\pi} > 1 \\ 0 & \hat{\pi} \leq 1 \end{cases} \Rightarrow \hat{Y} = 0$$

1 - Specificity = $P(\hat{Y}=1 | Y=0)$
= false positive rate
(want low)

C = "concordance index" = area under ROC curve
- estimates P (prediction & outcome are "concordant")

Higher $C \Rightarrow$
More "concordant"
= good
predictive power

larger $\hat{\pi} \rightarrow$ "larger" y

Framingham Example - Two predictors - Sbp, sex
Fitted additive model: quantitative | binary

$$\log\left(\frac{\hat{\pi}}{1-\hat{\pi}}\right) = -2.765 + 0.018 X_1 - 0.783 X_2$$

X_1 = systolic blood pressure (mmHg)

$X_2 = \begin{cases} 1 & \text{female} \\ 0 & \text{male} \end{cases}$

Interpretations in R code -

Males: $\log\left(\frac{\hat{\pi}}{1-\hat{\pi}}\right) = -2.765 + \underbrace{0.018 X_1}$

Females: $\log\left(\frac{\hat{\pi}}{1-\hat{\pi}}\right) = -2.765 - 0.783 + \underbrace{0.018 X_1}$

Females - Males : -0.783

Fitted interaction model:

$$\log\left(\frac{\hat{\pi}}{1-\hat{\pi}}\right) = -2.089 + 0.0128 X_1 - 1.867 X_2 + 0.0080 X_1 X_2$$

Males ($X_2=0$): $-2.089 + 0.0128 X_1$

Females ($X_2=1$): $(-2.089 - 1.867) + (0.0128 + 0.0080) X_1$

Females - Males : $-1.867 + 0.008 X_1$

$\log\left(\frac{\text{odds CHD Female} | X_1}{\text{odds CHD Male} | X_1}\right)$

How to interpret the interaction coef itself?

Multiplicative = effect on an effect
= change in odds ratio

effect of sex on
the effect of sbp on CHD

$$\frac{\text{odds of CHD sbp } x+1}{\text{odds of CHD sbp } x}$$

effect of sbp
on the effect of sex on CHD

$$\frac{\text{odds CHD Female}}{\text{odds CHD Males}}$$