

classification

Y categorical (not numeric)

ground truth: for given x ,

$$P_1(x_0) = P(Y=1 | X=x_0)$$

\vdots

$$P_k(x_0) = P(Y=k | X=x_0)$$

} randomness
Some x_0 can give different values of Y
irreducible error

want $\hat{f}(x_0) \rightarrow$ returns the category of Y most likely at x_0

measure error

misclassification rate: $\frac{1}{n} \sum_{i=1}^n 1\{Y_i \neq \hat{Y}_i\}$
between 0 and 1

accuracy = 1 - misclassification rate

Confusion matrix ($k=2$)

		true	
		1	0
pred	1	a	b
	0	c	d

$a = \# \text{ samples where } \hat{Y}_i = 1 \text{ and } Y_i = 1$
 $b = \dots \dots \dots \hat{Y}_i = 1 \quad Y_i = 0$
 $c = \dots \dots \dots \hat{Y}_i = 0 \quad Y_i = 1$
 $d = \dots \dots \dots \hat{Y}_i = 0 \quad Y_i = 0$

$$n = a + b + c + d$$

$$\text{misclassification rate} = \frac{b+c}{n}$$

$$\text{accuracy} = \frac{a+d}{n}$$

before, used CV to estimate MSPE

now use CV to estimate misclassification rate

$k=2$

compute $P(\hat{Y}=1|x)$, $P(\hat{Y}=0|x)$

predict $\hat{Y}=1$ if $P(\hat{Y}=1|x) > 0.5 \leftarrow \text{threshold}$

Why not consider different values of threshold?
might want to consider different threshold

\hookrightarrow can get smaller misclassification error with different threshold