

Instance  $\rightarrow$  (continuous risk) score  $\Delta$

Cost ratio  $c$   
 decision/classification  
 threshold  $t$

$$\text{cost}(t) = \text{FP}(t) + c \text{FN}(t)$$

$\hookrightarrow$  total cost of all decisions when  $t$  is applied

Minimize the expected cost by taking the derivative w.r.t.  $t$  and setting it to zero:

$f(s) \rightarrow$  density of cases w/ a given score

$E[Y|s] \rightarrow$  rate of  
 cancer cases (positives)  
 among the patients  
 for a given score

$$\frac{dE[\text{cost}(t)]}{dt} = \frac{dE[\text{FP}(t)]}{dt} + c \frac{dE[\text{FN}(t)]}{dt}$$

$$0 = -f(t^*)E[1-Y|s=t^*] + c f(t^*)E[Y|s=t^*]$$

$\hookrightarrow$  positive

$\uparrow$  threshold  $\therefore \downarrow$  FPs,  $\uparrow$  FNs  
 $\swarrow$   
 # cases at the threshold

$\searrow$   
 fraction of cases  
 that were negative

$$t^*_{\text{impl}} = E[Y|s=t^*] = \frac{1}{1+c}$$

$$\frac{1}{1+1} = \frac{1}{2} = 0.5$$

$\hookrightarrow$  implied threshold

score at which cases have a  $1/(1+c)$  probability of being positive.

$$\text{FNR} = \frac{\text{FN}}{\text{FN} + \text{TP}} = \frac{3}{3+4} = \frac{3}{7} \approx 0.43 \text{ male}$$

$$\frac{2}{2+7} = \frac{2}{9} \approx 0.22 \text{ Female}$$

group a  $\rightarrow$  publisher  
 group b  $\rightarrow$  consumer

$$\text{cost}_{ab} = \text{FP}_{ab}(t_{ab}) + c_{ab} \text{FN}_{ab}(t_{ab})$$

$$E[Y|s=t, \text{female}] = 0.25$$

$$c = \frac{1-0.25}{0.25} = 3$$

$$E[Y|s=t, \text{male}] = 0.16$$

$$c = (1-0.16)/0.16 \approx 5$$

$$c = \frac{1 - E[Y|s=t]}{E[Y|s=t]} \rightarrow \text{cost ratio}$$

$\hookrightarrow$  implied threshold

FN

Decision threshold  $t$

Score  $s$

Label  $\hat{y}$

$\hat{y}_i = 1_{\{s_i \geq t\}}$  ↑ indicator function

$t_{\text{implied}} = E[Y | s=t] \rightarrow$  prob. of the outcome corresponding to the decision threshold

(classifier)  
If calibrated,  $s = E[Y | s=s]$ , then the decision and implied thresholds are identical

$$E[Y | s=t] \approx E[Y | t_{\text{lower}} \leq s \leq t_{\text{upper}}], \quad t_{\text{lower}} \leq t < t_{\text{upper}}$$

↓ ↑ implied threshold estimation

$$Y_i = \beta_0 + \beta_1 (s_i - t) + \epsilon_i \quad \text{weighted linear regression}$$

$$w_i = \begin{cases} \left[ 1 - \left( \frac{|s_i - t|}{d} \right)^3 \right]^3 & s_i \in [t-d, t+d] \\ 0 & \text{otherwise} \end{cases}$$

half-width, symmetric region around the threshold

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

Fairness measurement dataset:

$(Y_{ij}, Y_i^*) \rightarrow$  expert-provided ground-truth for that post  
↓  
label provided by labeler  $j$  to post  $i$

Signal Detection Theory (SDT):

Parameters:  $\phi$  prevalence

$d'$  separation  $\rightarrow$  distance between the two means.

$$FPR(d', t) = 1 - \Phi(t)$$

$$FNR(d', t) = \Phi(t - d')$$

model  
parameters

↓  
cumulative dist. function for the standard normal

using the observed error rates:

$\phi \rightarrow$  prevalence

$$t = \Phi^{-1}(1 - FPR)$$

$$d' = t - \Phi^{-1}(FNR)$$

↪ inverse cumulative dist. function for the standard normal  
(Prevalence and Bayes' rule)

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$E[Y|s=t] = P(Y=1|s=t)$$

$$= \frac{P(Y=1) P(s=t|Y=1)}{P(Y=1) P(s=t|Y=1) + P(Y=0) P(s=t|Y=0)}$$

$$= \frac{\phi N(t-d')}{\phi N(t-d') + (1-\phi) N(t)}$$

Equal variance

$N(\cdot) \rightarrow$  density function  
of the standard normal  
dist.

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

implied  
threshold

$$C = \frac{1 - E[Y|s=t]}{E[Y|s=t]}$$

$$\begin{aligned} & \xrightarrow{(\div \phi N(t-d'))} = \frac{1}{1 + \frac{1-\phi}{\phi} \exp(-td' + d'^2/2)} \end{aligned}$$

$$C = \frac{1-\phi}{\phi} \exp(-td' + d'^2/2)$$

$$-\frac{1}{2}\left(\frac{x-t}{\sigma}\right)^2 - \left(-\frac{1}{2}\left(\frac{x-(t-d')}{\sigma}\right)^2\right)$$

$$-\frac{1}{2}(x-t)^2 + \frac{1}{2}(x-(t-d'))^2$$