

# BINARY CLASSIFICATION

STAT 432  
DALPIAZ

MODEL FIT TO ESTIMATION

$(x_i, y_i)$

FROM

VALIDATION → EVAL

OR

TEST → REPORT

MODEL FIT TO TRAIN

$y_i$   
↓  
ACTUAL  
O  
O  
O  
O  
O  
O  
O  
I  
I  
I  
I

$C(x_i)$   
↓  
PREDICTED  
O  
O  
O  
O  
I → "POSITIVE"  
I  
O → "NEGATIVE"  
I  
I  
I

# CONFUSION MATRIX

		<u>ACT</u>	
		I	O
<u>PRED</u>	I		
	O		

↑  
POSITIONS OF ACT/PRED  
AND O/I COULD CHANGE !!!

<u>ACTUAL</u>	<u>PREDICTED</u>	
0	0	→ TN
0	0	
0	0	
0	0	
0	1	→ FP
0	1	
1	0	→ FN
1	1	
1	1	
1	1	→ TP

		<u>ACT</u>	
		1	0
<u>PRED</u>	1	3	2
	0	1	4

TP (True Positive) points to the cell (1,1) containing 3.  
 FP (False Positive) points to the cell (1,0) containing 2.  
 FN (False Negative) points to the cell (0,1) containing 1.  
 TN (True Negative) points to the cell (0,0) containing 4.

P = 4 (Minority Class)  
 N = 6 (Majority Class)

MINORITY CLASS      MAJORITY CLASS

		<u>ACT</u>	
		1	0
<u>PRED</u>	1	3	2
	0	1	4

TP (True Positive) points to the cell (1,1) containing 3.  
 FP (False Positive) points to the cell (1,0) containing 2.  
 FN (False Negative) points to the cell (0,1) containing 1.  
 TN (True Negative) points to the cell (0,0) containing 4.  
 P = 4 (Total Positive Predictions) and N = 6 (Total Actual Positives) are indicated below the table.

\*  $\text{Acc} = \frac{\text{TP} + \text{TN}}{\text{P} + \text{N}} = \frac{3 + 4}{4 + 6} = 0.7$  (1 - MISCLASS)

AND MANY MORE... ( $F_1$ , MCC, ...)

SEE BOOK/WIKIPEDIA

$$P_{\text{REV}} = \frac{P}{P + N} = \frac{4}{4 + 6} = 0.4$$

$$S_{\text{ENS}} = \frac{\text{TP}}{P} = \frac{3}{4} = 0.75$$

$$S_{\text{PEC}} = \frac{\text{TN}}{N} = \frac{4}{6} = 0.66\bar{6}$$

$$P_{\text{PV}} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{3}{3 + 2} = 0.6$$

$$F_{\text{DR}} = \frac{\text{FP}}{\text{FP} + \text{TP}} = \frac{2}{2 + 3} = 0.4$$

No INFORMATION RATE → PROPORTION OF MAJORITY CLASS

$$NIR = \max \left\{ \overset{\text{PREV}}{\frac{P}{P+N}}, \frac{N}{P+N} \overset{1-\text{PREV}}{\quad} \right\}$$

IF  $ACC > NIR$  → REASONABLE CLASSIFIER

IF  $ACC < NIR$  → USELESS CLASSIFIER

$$0.7 > 0.6 \quad \checkmark$$

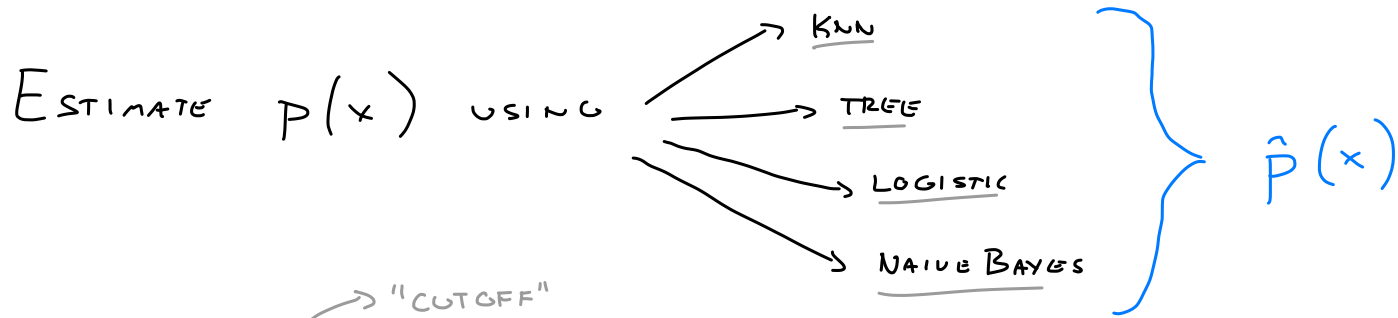
$$\underline{Y=0} \quad \text{or} \quad \underline{Y=1}$$

$$p(x) \triangleq P[Y=1 | X=x]$$

$$1 - p(x) = P[Y=0 | X=x]$$

$$\rightarrow P[Y=1 | X=x] \geq P[Y=0 | X=x]$$

$$C^B(x) = \begin{cases} 1 & p(x) \geq 0.5 \\ 0 & p(x) < 0.5 \end{cases}$$



SET  $0 \leq \alpha \leq 1$  "CUTOFF"

$$C_{\alpha}(x) = \begin{cases} 1 & \hat{P}(x) \geq \alpha \\ 0 & \hat{P}(x) < \alpha \end{cases}$$


ASSUMED  $\alpha=0.5$

AS  $\alpha \uparrow$  HARDER TO CLASSIFY AS  $Y=1$

$$\text{CLASSIFIER} = f(\hat{P}(x), \alpha)$$

$y_i$	$\hat{p}(x_i)$	$C_{0.0}(x_i)$	$C_{0.25}(x_i)$	$C_{0.5}(x_i)$	$C_{0.75}(x_i)$	$C_{1.0}(x_i)$
0	0.1	1	0	0	0	0
0	0.1	1	0	0	0	0
0	0.2	1	0	0	0	0
0	0.3	1	1	0	0	0
0	0.6	1	1	1	0	0
0	0.7	1	1	1	0	0
1	0.4	1	1	0	0	0
1	0.7	1	1	1	0	0
1	0.8	1	1	1	1	0
1	0.9	1	1	1	1	0

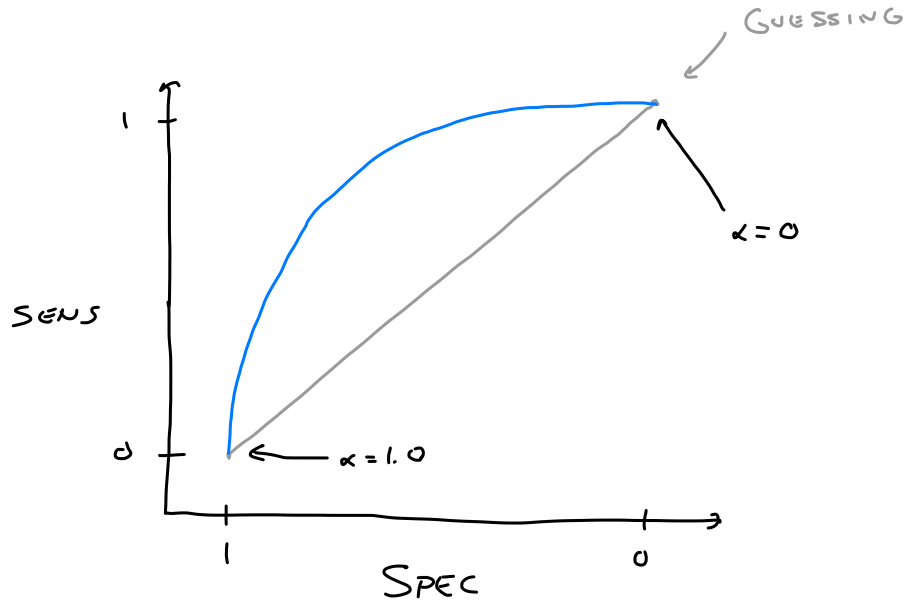
$TP/P = SENS =$	1.00	1.00	0.75	0.50	0.00
$TN/N = SPEC =$	0.00	0.50	0.66	1.00	1.00
$ACC =$	0.40	0.70	0.70	0.80	0.60


 EXPECTED TO BE BEST



EVALUATE  $\hat{p}(x)$  INSTEAD OF  $C(x)$ ? ROC CURVE!

INPUTS  
 $y$   
 $\hat{p}(x)$



AUC

- Bigger = Better
- 1 = PERFECT
- 0.50 = "WORST"