

# Classification Assessment

Let  $\vec{x} = \{x_1, x_2, \dots, x_d\}^T \in \mathbb{R}^d$   
 $\hat{y} \in \{c_1, c_2, \dots, c_k\}$

Given a classifier  $M: \mathbb{R}^d \rightarrow \{c_1, \dots, c_k\}$   
 created w/ training set

Assess  $M$  w/ test set, a set of pts w/ known labels

Let  $D$  be  $n$  pts in  $\mathbb{R}^d$

w/ known labels

let  $y_i$  be the known label of  $\vec{x}_i$   
 $\hat{y}_i = M(\vec{x}_i)$  is the predicted label

## Performance Measures

### Error rate

Let  $I$  be an indicator function

takes val 1 when arg is true false otherwise

$$\text{ErrorRate} = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i) \leftarrow \begin{array}{l} \text{- fraction of incorrect} \\ \text{predictions} \end{array}$$

- est of prob
- of misclassification
- low error rate
- is better

### Accuracy

$$\text{Accuracy} = 1 - \text{ErrorRate} \leftarrow \begin{array}{l} \text{fraction of correct predictions} \\ \text{est of the prob of correct} \\ \text{classification} \\ \text{high accuracy is better} \end{array}$$

What could be a problem w/ just error rate or accuracy?

## Contingency table based Measures

Let  $D$  as our test set

$$D_i \subset D \text{ w/ label } c_i$$

$$n_i = |D_i|$$

$$D = \{D_1, \dots, D_k\}$$

$R_j \subseteq D$  w/ predicted label  $c_i$

$$m_j = |R_j|$$

$$R = \{R_1, \dots, R_k\}$$

The **contingency table**  $N$  (aka Confusion matrix)

is  $k \times k$

$$\text{u/ } N(i,j) = n_{ij} = |R_i \cap D_j| = \left| \sum \vec{x}_k \in D \mid \hat{y}_k = c_i \text{ and } y_k = c_j \right|$$

# of pts w/  
predicted label  $c_i$

and true  
label  $c_j$

What is  $n_{ij}$ ?

# of pts w/  
predicted label  $c_i$

and true label  $c_i \Rightarrow$  # of times  
we got  $c_i$  correct ✓

What is  $\sum_{i=1, i \neq j}^k n_{ij}$  ?

# of times  $c_i$  wrong

# of pts

predicted label  $c_i$

actual label or  $c_j$

## Accuracy (Precision)

### Class specific accuracy (precision)

$$\text{acc}_i = \text{prec}_i = \frac{n_{ii}}{m_i} \quad \begin{matrix} \leftarrow \# \text{ of correct } c_i's \\ \leftarrow \# \text{ w/ label } c_i \end{matrix}$$

Higher is better

### Overall accuracy (precision)

$$\text{Accuracy} = \text{Precision} = \sum_{i=1}^k \frac{m_i}{n} \text{acc}_i = \frac{1}{n} \sum_{i=1}^k n_{ii}$$

Coverage / Recall

### Class specific coverage (recall)

$$\text{Coverage}_i = \text{recall}_i = \frac{n_{ii}}{n_i}$$

Higher is better

How can I get  $\text{recall}_i = 1$  w/ little effort?

$$\text{define } N(\vec{x}) = c_i \quad \forall \vec{x} \in \mathbb{R}^d$$

Problem: precision would be very low

## F-measure

class-specific F-measure : balance precision & recall

$$F_i = \frac{2}{\frac{1}{\text{rec}_i} + \frac{1}{\text{prec}_i}} = \frac{2 n_{ii}}{n_{ii} + m_i} \leftarrow \begin{matrix} \text{higher is} \\ \text{better!} \end{matrix}$$