

MODEL EVALUATION AND SELECTION

- Evaluation metrics: How can we measure accuracy?
Other metrics to consider?
- Use **validation test set** of class-labeled tuples instead of training set when assessing accuracy
- Some of the measures are:
 - Accuracy – suitable when class tuples are evenly distributed
 - Precision - suitable when class tuples are not evenly distributed
 - Recall - Sensitivity



CLASSIFIER EVALUATION METRICS: CONFUSION MATRIX

Confusion Matrix:

Actual class\Predicted class	Yes	No
Yes	True Positives (TP)	False Negatives (FN)
No	False Positives (FP)	True Negatives (TN)

Actual class\Predicted class	C_1	$\neg C_1$
C_1	True Positives (TP)	False Negatives (FN)
$\neg C_1$	False Positives (FP)	True Negatives (TN)

- Given m classes, an entry, $CM_{i,j}$ in a **confusion matrix** indicates # of tuples in class i that were labeled by the classifier as class j
- May have extra rows/columns to provide totals



CLASSIFIER EVALUATION METRICS: CONFUSION MATRIX

- True Positives
 - Positive tuples correctly classified as positive.
- True Negatives:
 - Negative tuples correctly classified as negative.
- False Positives:
 - Negative tuples incorrectly classified as positives.
- False Negatives:
 - Positive tuples incorrectly classified as negatives



CLASSIFIER EVALUATION METRICS: CONFUSION MATRIX

Confusion Matrix:

Actual class\Predicted class	C_1	$\neg C_1$
C_1	True Positives (TP)	False Negatives (FN)
$\neg C_1$	False Positives (FP)	True Negatives (TN)

Example of Confusion Matrix:

Actual class\Predicted class	buy_computer = yes	buy_computer = no	Total
buy_computer = yes	1	0	
buy_computer = no	1	998	
Total			1000

- Given m classes, an entry, $CM_{i,j}$ in a **confusion matrix** indicates # of tuples in class i that were labeled by the classifier as class j
- May have extra rows/columns to provide totals



ACCURACY, ERROR RATE, SENSITIVITY AND SPECIFICITY

A\P	C	¬C	
C	TP	FN	P
¬C	FP	TN	N
	P'	N'	All

- **Classifier Accuracy**, or recognition rate: percentage of test set tuples that are correctly classified

$$\text{Accuracy} = (TP + TN)/All$$

- **Error rate**: $1 - \text{accuracy}$, or
$$\text{Error rate} = (FP + FN)/All$$



CLASSIFIER EVALUATION METRICS: PRECISION AND RECALL, AND F-MEASURES

- **Precision:** exactness – what % of tuples that the classifier labeled as positive are actually positive

$$precision = \frac{TP}{TP + FP}$$

- **Recall:** completeness – what % of positive tuples did the classifier label as positive?
- Perfect score is 1.0

$$recall = \frac{TP}{TP + FN}$$



The 2-by-2 confusion matrix

		<i>gold standard labels</i>		
		gold positive	gold negative	
<i>system output labels</i>	system positive	true positive	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		recall = $\frac{tp}{tp+fn}$		accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

CONFUSION MATRIX FOR 3 CLASS CLASSIFICATION

		<i>gold labels</i>			
		urgent	normal	spam	
<i>system output</i>	urgent	8	10	1	precision_u = $\frac{8}{8+10+1}$
	normal	5	60	50	precision_n = $\frac{60}{5+60+50}$
	spam	3	30	200	precision_s = $\frac{200}{3+30+200}$
		recall_u = $\frac{8}{8+5+3}$	recall_n = $\frac{60}{10+60+30}$	recall_s = $\frac{200}{1+50+200}$	



Macroaveraging and Microaveraging

Class 1: Urgent

	true urgent	true not
system urgent	8	11
system not	8	340

$$\text{precision} = \frac{8}{8+11} = .42$$

Class 2: Normal

	true normal	true not
system normal	60	55
system not	40	212

$$\text{precision} = \frac{60}{60+55} = .52$$

Class 3: Spam

	true spam	true not
system spam	200	33
system not	51	83

$$\text{precision} = \frac{200}{200+33} = .86$$

Pooled

	true yes	true no
system yes	268	99
system no	99	635

$$\text{microaverage precision} = \frac{268}{268+99} = .73$$

$$\text{macroaverage precision} = \frac{.42+.52+.86}{3} = .60$$



- Macro-average gives equal weight to each class, regardless of the number of instances. Micro-averaging, on the other hand, aggregates the counts of true positives, false positives, and false negatives across all classes and then calculates the performance metric based on the total counts.



CLASSIFIER EVALUATION METRICS: PRECISION AND RECALL, AND F-MEASURES

- The F-score is a way of combining the precision and recall of the model, and it is defined as the harmonic mean of the model's precision and recall.
- Inverse relationship between precision & recall
- **F measure (F_1 or F-score)**: harmonic mean of precision and recall,

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

- F_β : weighted measure of precision and recall
 - assigns β times as much weight to recall as to precision

$$F_\beta = \frac{(1 + \beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$$

CLASSIFIER EVALUATION METRICS: EXAMPLE

Actual Class\Predicted class	cancer = yes	cancer = no	Total	Recognition(%)
cancer = yes	90	210	300	30.00 (<i>sensitivity</i>)
cancer = no	140	9560	9700	98.56 (<i>specificity</i>)
Total	230	9770	10000	96.40 (<i>accuracy</i>)

- $Precision = 90/230 = 39.13\%$

$$Recall = 90/300 = 30.00\%$$

