# Evaluation and Improvement of Classification Quality

# Outline

❏ Model Evaluation and Selection

❏ Techniques to Improve Classification Accuracy: Ensemble Methods

❏ Multiclass Classification and Weak Supervision

# Model Evaluation and Selection

# Model Evaluation and Selection

❑ How to evaluate the quality of a classifier

   ❑ A typical measure: Accuracy

   ❑ Other metrics to consider?

❑ How to assess the classification quality

   ❑ Use (independent) **validation test set** instead of training set when assessing accuracy

❑ Methods for estimating a classifier's accuracy

   ❑ Holdout method

   ❑ Cross-validation

   ❑ Bootstrap

❑ Comparing classifiers using ROC Curves

# 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)** |

❑ In a confusion matrix with $m$ classes, $CM_{i,j}$ indicates # of tuples in class $i$ that were labeled by the classifier as class $j$

  ❑ May have extra rows/columns to provide totals

❑ An example of Confusion Matrix:

| Actual class\Predicted class | buy_computer = yes | buy_computer = no | Total |
|---|---|---|---|
| buy_computer = yes | **6954** | **46** | 7000 |
| buy_computer = no | **412** | **2588** | 3000 |
| Total | 7366 | 2634 | 10000 |

# Classifier Evaluation Metrics: 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

  **Accuracy = (TP + TN)/All**
- **Error rate:** *1 – accuracy*, or

  **Error rate = (FP + FN)/All**

- ❑ **Class imbalance problem**
  - ❑ One class may be *rare*
    - ❑ E.g., fraud, or HIV-positive
  - ❑ Significant *majority of the negative class* and minority of the positive class
- ❑ Measures handle the class imbalance problem
  - ❑ **Sensitivity** (recall): True positive recognition rate
    - ❑ **Sensitivity = TP/P**
  - ❑ **Specificity**: True negative recognition rate
    - ❑ **Specificity = TN/N**

# Classifier Evaluation Metrics: Precision and Recall, and F-measures

❑ **Precision** (Exactness): What percentage of tuples labeled as positive is actually positive?

$$P = \text{Precision} = \frac{TP}{TP + FP}$$

❑ **Recall** (Completeness): What percentage of positive tuples are labeled as positive?

❑ Range: [0, 1]

$$R = \text{Recall} = \frac{TP}{TP + FN}$$

❑ The "inverse" relationship between precision & recall

❑ *F* **measure (**or *F-score*)**: Harmonic mean of precision and recall

❑ In general, it is the weighted measure of precision & recall

$$F_{\beta} = \frac{1}{\alpha \cdot \frac{1}{P} + (1 - \alpha) \cdot \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

Assigning β times as much weight to recall as to precision

❑ *F1-measure (balanced F-measure)*

❑ That is, when β = 1,

$$F_1 = \frac{2PR}{P + R}$$

# Classifier Evaluation Metrics: Example

❑ Use the same confusion matrix, calculate the measure just introduced

| Actual Class\Predicted class | cancer = yes | cancer = no | Total | Recognition(%) |
|---|---|---|---|---|
| cancer = yes | **90** | **210** | 300 | 30.00 (*sensitivity*) |
| cancer = no | **140** | **9560** | 9700 | 98.56 (*specificity*) |
| Total | 230 | 9770 | 10000 | 96.50 (*accuracy*) |

❑ Sensitivity = TP/P = 90/300 = 30%

❑ Specificity = TN/N = 9560/9700 = 98.56%

❑ Accuracy = (TP + TN)/All = (90+9560)/10000 = 96.50%

❑ Error rate = (FP + FN)/All = (140 + 210)/10000 = 3.50%

❑ Precision = TP/(TP + FP) = 90/(90 + 140) = 90/230 = 39.13%

❑ Recall = TP/ (TP + FN) = 90/(90 + 210) = 90/300 = 30.00%

❑ F1 = 2 P × R /(P + R) = 2 × 39.13% × 30.00%/(39.13% + 30%) = 33.96%

# Classifier Evaluation: Holdout & Cross-Validation

❑ **Holdout method**

  ❑ The given data set is randomly partitioned into two independent sets

    ❑ Training set (e.g., 2/3) for model construction

    ❑ Test set (e.g., 1/3) for accuracy estimation

  ❑ Repeated random sub-sampling validation: A variation of holdout

    ❑ Repeat holdout k times, accuracy = average of the accuracies obtained

❑ **Cross-validation** (*k*-fold, where k = 10 is most popular)

  ❑ Randomly partition the data into *k mutually exclusive* subsets, each approximately equal size

  ❑ At *i*-th iteration, use $D_i$ as test set and others as training set

  ❑ <u>Leave-one-out</u>: *k* folds where *k* = # of tuples, for small sized data

  ❑ <u>**Stratified cross-validation**</u>:  Folds are stratified so that class distribution in each fold is approximately the same as that in the initial data
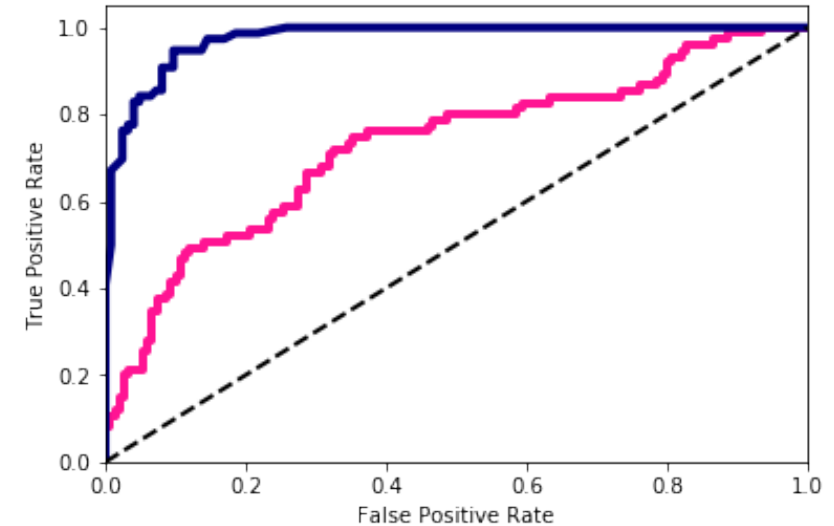
# Classifier Evaluation: Bootstrap

❑ **Bootstrap**

   ❑ Works well with small data sets

   ❑ Samples the given training tuples uniformly *with replacement*

   ❑ Each time a tuple is selected, it is equally likely to be selected again and re-added to the training set

❑ Several bootstrap methods, and a common one is **.632 bootstrap**

   ❑ A data set with *d* tuples is sampled *d* times, with replacement, resulting in a training set of *d* samples. The data tuples that did not make it into the training set end up forming the test set. About 63.2% of the original data end up in the bootstrap, and the remaining 36.8% form the test set (since $(1 - 1/d)^d \approx e^{-1} = 0.368$)

   ❑ Repeating the sampling procedure *k* times, the overall accuracy of the model:

$$Acc(M) = \frac{1}{k} \sum_{i=1}^{k} (0.632 \times Acc(M_i)_{test\_set} + 0.368 \times Acc(M_i)_{train\_set})$$

# Model Selection: ROC Curves

- ❑ **ROC** (Receiver Operating Characteristics) curve: for visual comparison of classification models

- ❑ Originated from signal detection theory

- ❑ Shows the trade-off between the true positive rate and the false positive rate

- ❑ The area under the ROC curve (**AUC**: Area Under Curve) is a measure of the accuracy of the model

- ❑ Rank the test tuples in decreasing order: The one that is most likely to belong to the positive class appears at the top of the list

- ❑ The closer to the diagonal line (i.e., the closer the area is to 0.5), the less accurate is the model



- ❑ The vertical axis represents the true positive rate
- ❑ The horizontal axis represents the false positive rate
- ❑ The plot also shows a diagonal line
- ❑ A model with perfect accuracy will have an area of 1.0

# Issues Affecting Model Selection

- **Accuracy**
  - Classifier accuracy: Predicting class label
- **Speed**
  - Time to construct the model (training time)
  - Time to use the model (classification/prediction time)
- **Robustness**: Handling noise and missing values
- **Scalability**: Efficiency in disk-resident databases
- **Interpretability**
  - Understanding and insight provided by the model
- Other measures
  - E.g., goodness of rules, such as decision tree size or compactness of classification rules