# Knowledge Discovery and Data Mining

Unit # 5

# Acknowledgement

- Most of the slides in this presentation are taken from course slides provided by
  - Han and Kimber (Data Mining Concepts and Techniques) and
  - Tan, Steinbach and Kumar (Introduction to Data Mining)

# Accuracy or Error Rates

- Partition: Training-and-testing
  - use two independent data sets, e.g., training set (2/3), test set(1/3)
  - used for data set with large number of examples

# Metrics for Performance Evaluation

- Focus on the predictive capability of a model
  - Rather than how fast it takes to classify or build models, scalability, etc.
- Confusion Matrix:

| | PREDICTED CLASS | |
|---|---|---|
| | Class=Yes | Class=No |
| ACTUAL CLASS Class=Yes | a | b |
| ACTUAL CLASS Class=No | c | d |

a: TP (true positive)
b: FN (false negative)
c: FP (false positive)
d: TN (true negative)

## Metrics for Performance Evaluation…

| | PREDICTED CLASS | | |
|---|---|---|---|
| ACTUAL CLASS | | Class=Yes | Class=No |
| | Class=Yes | a (TP) | b (FN) |
| | Class=No | c (FP) | d (TN) |

- Most widely-used metric:

$$\text{Accuracy} = \frac{a+d}{a+b+c+d} = \frac{TP+TN}{TP+TN+FP+FN}$$

Sajjad Haider                                Fall 2011                                5

## Limitation of Accuracy

- Consider a 2-class problem
  - Number of Class 0 examples = 9990
  - Number of Class 1 examples = 10

- If model predicts everything to be class 0, accuracy is 9990/10000 = 99.9 %
  - Accuracy is misleading because model does not detect any class 1 example

Sajjad Haider                                Fall 2011                                6

# Cost Matrix

|  | PREDICTED CLASS | | |
| --- | --- | --- | --- |
|  | C(i\|j) | **Class=Yes** | **Class=No** |
| **ACTUAL CLASS** | **Class=Yes** | C(Yes\|Yes) | C(No\|Yes) |
| | **Class=No** | C(Yes\|No) | C(No\|No) |

C(i|j): Cost of misclassifying class j example as class i

# Cost Matrix (Cont'd)

|  | PREDICTED CLASS | |
| --- | --- | --- |
|  | True | False |
| **ACTUAL CLASS** True | 10 | 5 |
| False | 1 | 14 |

|  | PREDICTED CLASS | |
| --- | --- | --- |
|  | True | False |
| **ACTUAL CLASS** True | 10 | 3 |
| False | 3 | 14 |

|  | PREDICTED CLASS | |
| --- | --- | --- |
|  | True | False |
| **ACTUAL CLASS** True | 10 | 6 |
| False | 0 | 14 |

All three confusion matrices have the same accuracy value, i.e., **24 / 30**

What if the cost of misclassification is not the same for both type of errors?

4

# Cost Matrix (Cont'd)

| | PREDICTED CLASS | |
|---|---|---|
| | True | False |
| ACTUAL CLASS — True | 10 | 5x5 |
| ACTUAL CLASS — False | 1 | 14 |

| | PREDICTED CLASS | |
|---|---|---|
| | True | False |
| ACTUAL CLASS — True | 10 | 3x5 |
| ACTUAL CLASS — False | 3 | 14 |

| | PREDICTED CLASS | |
|---|---|---|
| | True | False |
| ACTUAL CLASS — True | 10 | 6x5 |
| ACTUAL CLASS — False | 0 | 14 |

Suppose the cost of misclassifying True as False is 5 while the cost of misclassifying False as True is 1.

Accuracy values are:
**24/50, 24/42, 24/54**

# Cost Matrix (Cont'd)

| | PREDICTED CLASS | |
|---|---|---|
| | True | False |
| ACTUAL CLASS — True | 10 | 5x4 |
| ACTUAL CLASS — False | 1 | 14 |

| | PREDICTED CLASS | |
|---|---|---|
| | True | False |
| ACTUAL CLASS — True | 10 | 3x4 |
| ACTUAL CLASS — False | 3 | 14 |

| | PREDICTED CLASS | |
|---|---|---|
| | True | False |
| ACTUAL CLASS — True | 10 | 6x4 |
| ACTUAL CLASS — False | 0 | 14 |

Suppose the cost of misclassifying True as False is **4** while the cost of misclassifying False as True is 1.

Accuracy values are:
**24/45, 24/39, 24/48**

# Cost-Sensitive Measures

$$\text{Precision (p)} = \frac{a}{a+c}$$

$$\text{Recall (r)} = \frac{a}{a+b}$$

$$\text{F-measure (F)} = \frac{2rp}{r+p} = \frac{2a}{2a+b+c}$$

- Precision is biased towards C(Yes|Yes) & C(Yes|No)
- Recall is biased towards C(Yes|Yes) & C(No|Yes)
- F-measure is biased towards all except C(No|No)

$$\text{Weighted Accuracy} = \frac{w_1 a + w_4 d}{w_1 a + w_2 b + w_3 c + w_4 d}$$

Sajjad Haider                               Fall 2011                                    11

# Recall and Precision

| Actual | Prediction |
|--------|------------|
| T | T |
| T | F |
| F | T |
| F | F |
| F | T |
| T | T |
| T | T |
| T | F |
| F | T |
| T | T |

Sajjad Haider                               Fall 2011                                    12

# Recall and Precision

| Actual | Prediction |
|:------:|:----------:|
| T | T |
| T | F |
| F | T |
| F | F |
| F | T |
| T | T |
| T | T |
| T | F |
| F | T |
| T | T |

- Recall = 4 / 6

# Recall and Precision

| Actual | Prediction |
|:------:|:----------:|
| T | T |
| T | F |
| F | T |
| F | F |
| F | T |
| T | T |
| T | T |
| T | F |
| F | T |
| T | T |

- Recall = 4 / 6
- Precision = 4 / 7
- F-Measure = 8 / 13

# Terminology

- True Positive: The number of positive examples correctly predicted by the classification model.
- False Negative: The number of positive examples wrongly predicted as negative by the classification model.
- False Positive: The number of negative examples wrongly predicted as positive by the classification model.
- True Negative: The number of negative examples correctly predicted by the classification model.

# Terminology (Cont'd)

- The true positive rate (TPR) or sensitivity is defined as TPR = TP / (TP + FN).
- The true negative rate (TNR) or specificity is defined as TNR = TN / (TN + FP).
- The false positive rate (FPR) is defined as FPR = FP / (TN + FP).
- The false negative rate (FNR) is defined as FNR = FN / (TP + FN).

# ROC (Receiver Operating Characteristic)

- Developed in 1950s for signal detection theory to analyze noisy signals
  – Characterize the trade-off between positive hits and false alarms
- ROC curve plots TPR (on the y-axis) against FPR (on the x-axis)
- Remember that TPR represents "sensitivity" while FPR represents "100 – specificity".
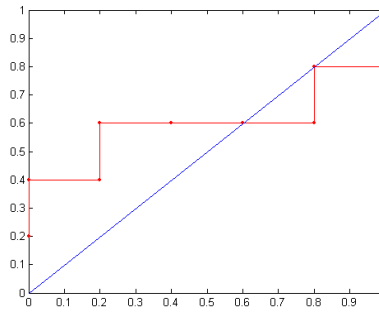
# How to Construct an ROC curve

| Instance | P(+|A) | True Class |
|----------|--------|------------|
| 1 | 0.95 | + |
| 2 | 0.93 | + |
| 3 | 0.87 | - |
| 4 | 0.85 | - |
| 5 | 0.85 | - |
| 6 | 0.85 | + |
| 7 | 0.76 | - |
| 8 | 0.53 | + |
| 9 | 0.43 | - |
| 10 | 0.25 | + |

- Use classifier that produces posterior probability for each test instance P(+|A)

- Sort the instances according to P(+|A) in decreasing order

- Apply threshold at each unique value of P(+|A)

- Count the number of TP, FP, TN, FN at each threshold

- TP rate, TPR = TP/(TP+FN)

- FP rate, FPR = FP/(FP + TN)

# How to construct an ROC curve

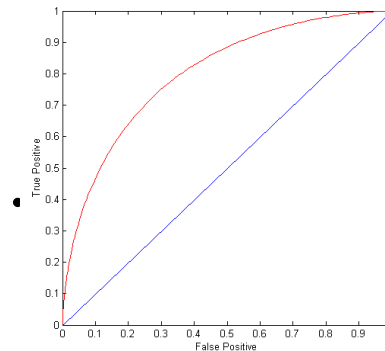| Class | + | - | + | - | - | - | + | - | + | + | |
|-------|---|---|---|---|---|---|---|---|---|---|---|
| Threshold >= | 0.25 | 0.43 | 0.53 | 0.76 | 0.85 | 0.85 | 0.85 | 0.87 | 0.93 | 0.95 | 1.00 |
| TP | 5 | 4 | 4 | 3 | 3 | 3 | 3 | 2 | 2 | 1 | 0 |
| FP | 5 | 5 | 4 | 4 | 3 | 2 | 1 | 1 | 0 | 0 | 0 |
| TN | 0 | 0 | 1 | 1 | 2 | 3 | 4 | 4 | 5 | 5 | 5 |
| FN | 0 | 1 | 1 | 2 | 2 | 2 | 2 | 3 | 3 | 4 | 5 |
| TPR | 1 | 0.8 | 0.8 | 0.6 | 0.6 | 0.6 | 0.6 | 0.4 | 0.4 | 0.2 | 0 |
| FPR | 1 | 1 | 0.8 | 0.8 | 0.6 | 0.4 | 0.2 | 0.2 | 0 | 0 | 0 |

ROC Curve:

# ROC Curve

- It shows the tradeoff between sensitivity and specificity (any increase in sensitivity will be accompanied by a decrease in specificity).
- The closer the curve follows the left-hand border and then the top border of the ROC space, the more accurate the test.
- The closer the curve comes to the 45-degree diagonal of the ROC space, the less accurate the test.
- The area under the curve is a measure of text accuracy.

# Lift and Gain Charts

- Very commonly used in the marketing research.
- **Lift** is a measure of the effectiveness of a predictive model calculated as the ratio between the results obtained with and without the predictive model.
- A lift chart consists of a lift curve and a baseline
- The greater the area between the lift curve and the baseline, the better the model

# Example

http://www2.cs.uregina.ca/~dbd/cs831/notes/lift_chart/lift_chart.html

- Using the response model P($x$)=100-AGE($x$) for customer $x$ and the data table, construct the cumulative gains and lift charts. Ties in ranking should be arbitrarily broken by assigning a higher rank to who appears first in the table.

| Customer Name | Height | Age | Actual Response |
|---|---|---|---|
| Alan | 70 | 39 | N |
| Bob | 72 | 21 | Y |
| Jessica | 65 | 25 | Y |
| Elizabeth | 62 | 30 | Y |
| Hilary | 67 | 19 | Y |
| Fred | 69 | 48 | N |
| Alex | 65 | 12 | Y |
| Margot | 63 | 51 | N |
| Sean | 71 | 65 | Y |
| Chris | 73 | 42 | N |
| Philip | 75 | 20 | Y |
| Catherine | 70 | 23 | N |
| Amy | 69 | 13 | N |
| Erin | 68 | 35 | Y |
| Trent | 72 | 55 | N |
| Preston | 68 | 25 | N |
| John | 64 | 76 | N |
| Nancy | 64 | 24 | Y |
| Kim | 72 | 31 | N |
| Laura | 62 | 29 | Y |

# Example: Steps 1 & 2

1. Calculate P(*x*) for each person *x*

2. Order the people according to rank P(*x*)

| Customer Name | P(x) | Actual Response |
|---|---|---|
| Alex | 88 | Y |
| Amy | 87 | N |
| Hilary | 81 | Y |
| Philip | 80 | Y |
| Bob | 79 | Y |
| Catherine | 77 | N |
| Nancy | 76 | Y |
| Jessica | 75 | Y |
| Preston | 75 | N |
| Laura | 71 | Y |
| Elizabeth | 70 | Y |
| Kim | 69 | N |
| Erin | 65 | Y |
| Alan | 61 | N |
| Chris | 58 | N |
| Fred | 52 | N |
| Margot | 49 | N |
| Trent | 45 | N |
| Sean | 35 | Y |
| John | 24 | N |

Sajjad Haider                    Fall 2011

# Example: Step 3

- Calculate the percentage of total responses for each cutoff point
  - Response Rate = Number of Responses / Total Number of Responses (10)

| Total Customers Contacted | Number of Responses | Response Rate |
|---|---|---|
| 2 | 1 | 10% |
| 4 | 3 | 30% |
| 6 | 4 | 40% |
| 8 | 6 | 60% |
| 10 | 7 | 70% |
| 12 | 8 | 80% |
| 14 | 9 | 90% |
| 16 | 9 | 90% |
| 18 | 9 | 90% |
| 20 | 10 | 100% |

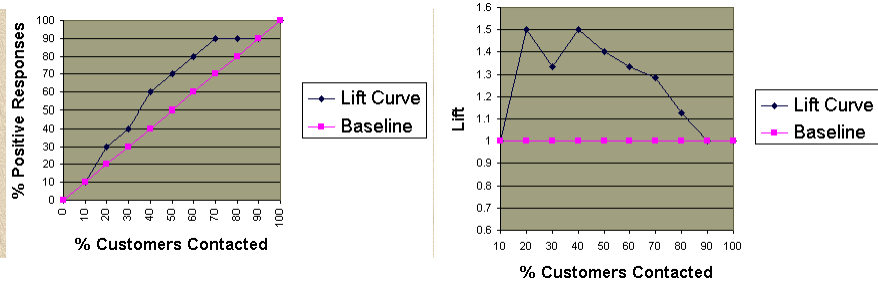Sajjad Haider                    Fall 2011                    24

# Example: Gains and Lift Charts

- To plot life chart, calculate the points on the lift curve by determining the ratio between the result predicted by our model and the result using no model.

# Practice Exercises

- Draw gains and lift charts of data set given on slide 18 .
- Draw ROC for data set given in slide 23.
- Apply supervised and unsupervised discretization techniques on "Attribute 2" of data set given on slide 35 of Unit # 3.
- Form decision trees using Entropy, Gini and Gain_ratio as splitting criteria on the above data set (after discretization).