# COSC 522 – Machine Learning

# Lecture 9 – Classifier Fusion

Hairong Qi, Gonzalez Family Professor
Electrical Engineering and Computer Science
University of Tennessee, Knoxville
https://www.eecs.utk.edu/people/hairong-qi/
Email: hqi@utk.edu

Course Website: http://web.eecs.utk.edu/~hqi/cosc522/

# Roadmap

- Supervised learning
  - Maximum Posterior Probability (MPP): For a given x, if $P(w_1|x) > P(w_2|x)$, then x belongs to class 1, otherwise 2.
    - Parametric Learning
      - Case 1: Minimum Euclidean Distance (Linear Machine), $\Sigma_i = \sigma^2 I$
      - Case 2: Minimum Mahalanobis Distance (Linear Machine), $\Sigma_i = \Sigma$
      - Case 3: Quadratic classifier, $\Sigma_i$ = arbitrary
      - Estimate Gaussian parameters using MLE
    - Nonparametric Learning
      - Parzon window (fixed window size)
      - K-Nearest Neighbor (variable window size)
- Unsupervised learning
  - Non-probabilistic approaches
    - kmeans, wta
  - Hierarchical approaches
    - Agglomerative clustering

- Supporting preprocessing techniques
  - Dimensionality Reduction
    - Supervised linear (FLD)
    - Unsupervised linear (PCA)
    - Unsupervised nonlinear (t-SNE)
- Supporting postprocessing techniques
  - Classifier Fusion

# Questions

- Rationale with fusion?
- Different flavors of fusion?
- The fusion hierarchy
- What is the cost function for Naïve Bayes?
- What is the procedure for Naïve Bayes?
- What is the limitation of Naïve Bayes?
- What is the procedure of Behavior-Knowledge-Space (BKS)?
- How does it resolve issues with NB?
- What is Boosting and what is its difference to committee-based fusion approaches?
- What is AdaBoost?

THE UNIVERSITY OF TENNESSEE KNOXVILLE
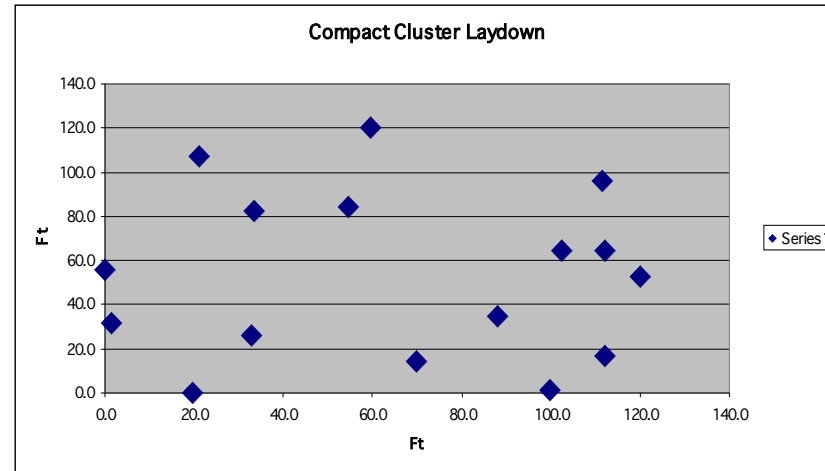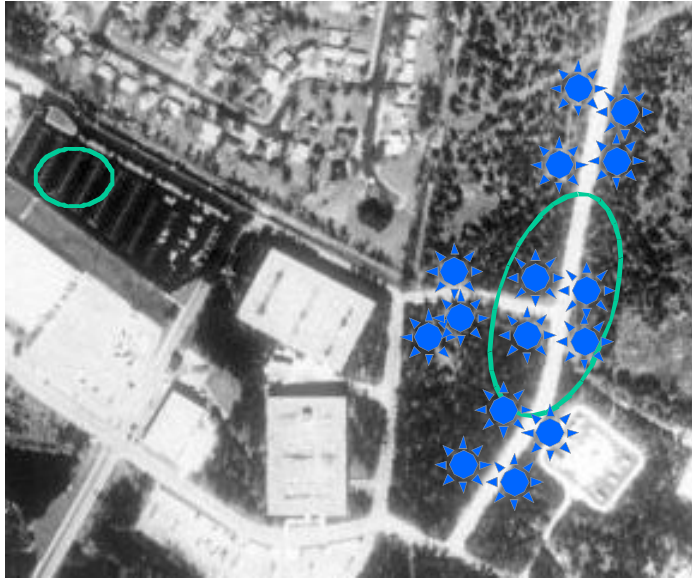
# Motivation

**Three heads are better than one.**

- Combining classifiers to achieve higher accuracy
  - Combination of multiple classifiers
  - Classifier fusion
  - Mixture of experts
  - Committees of neural networks
  - Consensus aggregation
  - …
- Reference:
  - L. I. Kuncheva, J. C. Bezdek, R. P. W. Duin, "Decision templates for multiple classifier fusion: an experimental comparison," *Pattern Recognition*, 34: 299-314, 2001.
  - Y. S. Huang and C. Y. Suen, "A method of combining multiple experts for the recognition of unconstrained handwritten numerals," IEEE Trans. Pattern Anal. Mach. Intell., vol. 17, no. 1, pp. 90–94, Jan. 1995.

# Popular Approaches

◆ Data-based fusion (early fusion)

◆ Feature-based fusion (middle fusion)

◆ Decision-based fusion (late fusion)

◆ Approaches

    ◆ Committee-based

        ◆ Majority voting

        ◆ Bootstrap aggregation (Bagging) [Breiman, 1996]

    ◆ Baysian-based

        ◆ Naïve Bayes combination (NB)

        ◆ Behavior-knowledge space (BKS) [Huang and Suen, 1995]

    ◆ Boosting

        ◆ Adaptive boosting (AdaBoost) [Freund and Schapire, 1996]

    ◆ Interval-based integration

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Application Example – Civilian Target Recognition

Compact Cluster Laydown

Ford 250     Harley Motocycle     Ford 350     Suzuki Vitara

THE UNIVERSITY OF TENNESSEE KNOXVILLE

# Consensus Patterns

- Unanimity (100%)
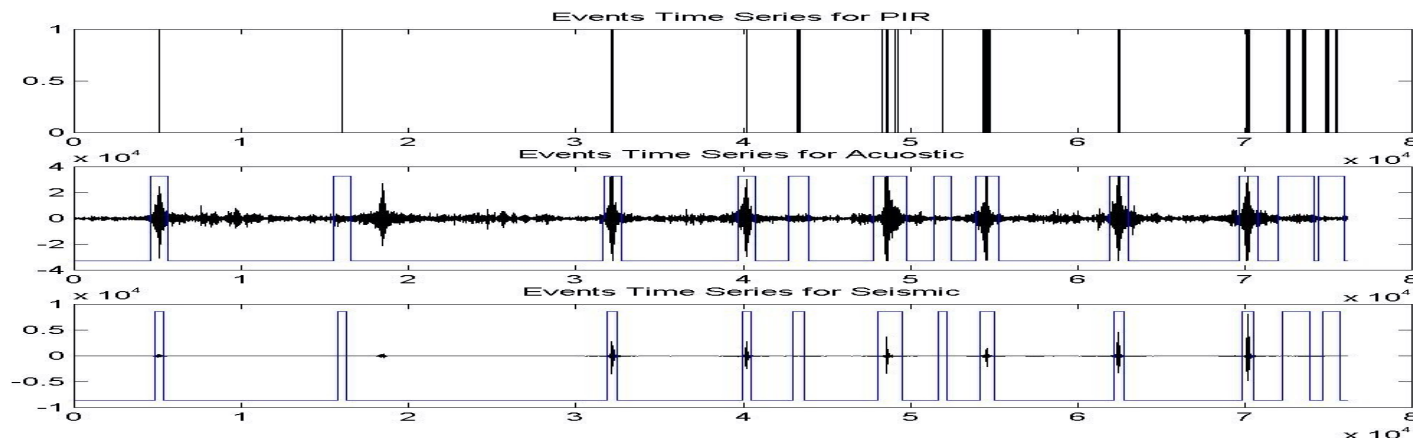- Simple majority (50%+1)
- Plurality (most votes)

# Example of Majority Voting - Temporal Fusion

◆ Fuse all the 1-sec sub-interval local processing results corresponding to the same event (usually lasts about 10-sec)

◆ Majority voting

$$\overline{\varphi_i^j} = \arg\max \omega_c, \quad c \in [1, C]$$

number of local
output c occurrence

number of possible local
processing results

# Questions

- Rationale with fusion?
- Different flavors of fusion?
- The fusion hierarchy
- What is the cost function for Naïve Bayes?
- What is the procedure for Naïve Bayes?
- What is the limitation of Naïve Bayes?
- What is the procedure of Behavior-Knowledge-Space (BKS)?
- How does it resolve issues with NB?
- What is Boosting and what is its difference to committee-based fusion approaches?
- What is AdaBoost?

THE UNIVERSITY OF TENNESSEE KNOXVILLE

# Naïve Bayes (the independence assumption)

The real class is DW, the classifier says it's HMV

Confusion matrix

s

| C1 | AAV | DW | HMV |
|---|---|---|---|
| AAV | 894 | 329 | 143 |
| DW | 99 | 411 | 274 |
| HMV | 98 | 42 | 713 |

k

| C2 | AAV | DW | HMV |
|---|---|---|---|
| AAV | 1304 | 156 | 77 |
| DW | 114 | 437 | 83 |
| HMV | 13 | 107 | 450 |

i = 1, 2 (classifiers)

| L1 | AAV | DW | HMV |
|---|---|---|---|
| AAV | | | |
| DW | | | |
| HMV | | | |

| L2 | AAV | DW | HMV |
|---|---|---|---|
| AAV | | | |
| DW | | | |
| HMV | | | |

Probability that the true class is k given that $C_i$ assigns it to s

Probability multiplication

THE UNIVERSITY OF TENNESSEE KNOXVILLE

10

# NB – Derivation

- Assume the classifiers are mutually independent
- Bayes combination - Naïve Bayes, simple Bayes, idiot's Bayes
- Assume
  - L classifiers, i=1,..,L
  - c classes, k=1,…,c
  - $s_i$: class label given by the i[th] classifier, i=1,…,L, s={$s_1$,…,$s_L$}

$$P(\omega_k|\mathbf{s}) = \frac{p(\mathbf{s}|\omega_k)P(\omega_k)}{p(\mathbf{s})} = \frac{P(\omega_k)\prod_{i=1}^{L}p(s_i|\omega_k)}{p(\mathbf{s})}$$

$$P(\omega_k) = N_k/N$$

$$p(s_i|\omega_k) = cm_{k,s_i}/N_k$$

$$P(\omega_k|\mathbf{s}) \approx \frac{1}{N_k^{L-1}}\prod_{i=1}^{L}cm_{k,s_i}$$

# BKS

- Majority voting won't work
- Behavior-Knowledge Space algorithm (Huang&Suen)

Assumption:
- 2 classifiers
- 3 classes
- 100 samples in the training set
Then:
- 9 possible classification combinations

| $c_1, c_2$ | samples from each class | fused result |
|---|---|---|
| 1,1 | 10/3/3 | 1 |
| 1,2 | 3/0/6 | 3 |
| 1,3 | 5/4/5 | 1,3 |
| ... | | |
| 3,3 | 0/0/6 | 3 |

# Questions

- Rationale with fusion?
- Different flavors of fusion?
- The fusion hierarchy
- What is the cost function for Naïve Bayes?
- What is the procedure for Naïve Bayes?
- What is the limitation of Naïve Bayes?
- What is the procedure of Behavior-Knowledge-Space (BKS)?
- How does it resolve issues with NB?
- What is Boosting and what is its difference to committee-based fusion approaches?
- What is AdaBoost?

# Boosting

- Base classifiers are trained in sequence!
- Base classifiers as weak learners
- Weighted majority voting to combine classifiers



$$Y_M(\mathbf{x}) = \text{sign}\left(\sum_m^M \alpha_m y_m(\mathbf{x})\right)$$

# AdaBoost

- Step 1: Initialize the data weighting coefficients $\{w_n\}$ by setting $w_n^{(1)} = 1/N$, where $N$ is the # of samples

- Step 2: for each classifier $y_m(\mathbf{x})$

  - (a) Fit a classifier $y_m(\mathbf{x})$ to the training data by minimizing the weighted error function

$$J_m = \sum_{n=1}^{N} w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)$$

  - (b) Evaluate the quantities

$$\epsilon_m = \frac{\sum_{n=1}^{N} w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^{N} w_n^{(m)}} \qquad \alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}$$

  - (c) Update the data weighting coefficients

$$w_n^{(m+1)} = w_n^{(m)} \exp \left\{ \alpha_m I(y_m(\mathbf{x}_n) \neq t_n) \right\}$$

- Step 3: Make predictions using the final model

$$Y_M(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^{M} \alpha_m y_m(\mathbf{x}) \right)$$

**Figure 14.2** Illustration of boosting in which the base learners consist of simple thresholds applied to one or other of the axes. Each figure shows the number $m$ of base learners trained so far, along with the decision boundary of the most recent base learner (dashed black line) and the combined decision boundary of the ensemble (solid green line). Each data point is depicted by a circle whose radius indicates the weight assigned to that data point when training the most recently added base learner. Thus, for instance, we see that points that are misclassified by the $m = 1$ base learner are given greater weight when training the $m = 2$ base learner.

# Questions

- Rationale with fusion?
- Different flavors of fusion?
- The fusion hierarchy
- What is the cost function for Naïve Bayes?
- What is the procedure for Naïve Bayes?
- What is the limitation of Naïve Bayes?
- What is the procedure of Behavior-Knowledge-Space (BKS)?
- How does it resolve issues with NB?
- What is Boosting and what is its difference to committee-based fusion approaches?
- What is AdaBoost?
- Interval-based fusion

# Value-based vs. Interval-based Fusion

- Interval-based fusion can provide fault tolerance
- Interval integration – overlap function
  - Assume each sensor in a cluster measures the same parameters, the integration algorithm is to construct a simple function (overlap function) from the outputs of the sensors in a cluster and can resolve it at different resolutions as required



Crest: the highest, widest peak of the overlap function

18

# A Variant of kNN

- Generation of local confidence ranges (For example, at each node i, use kNN for each $k \in \{5, \ldots, 15\}$)

| | Class 1 | Class 2 | … | Class n | |
|---|---|---|---|---|---|
| **k=5** | 3/5 | 2/5 | … | 0 | ← confidence level |
| **k=6** | 2/6 | 3/6 | … | 1/6 | |
| **…** | … | … | … | … | |
| **k=15** | 10/15 | 4/15 | … | 1/15 | |

**{2/6, 10/15}**  **{4/15, 3/6}**  …  **{0, 1/6}** ← confidence range

smallest  largest in this column

- Apply the integration algorithm on the confidence ranges generated from each node to construct an overlapping function

# Example of Interval-based Fusion

|         | stop 1 | | stop 2 | | stop 3 | | stop 4 | |
|---------|------|-------|------|-------|------|-------|------|-------|
|         | c    | acc   | c    | acc   | c    | acc   | c    | acc   |
| class 1 | 1    | 0.2   | 0.5  | 0.125 | 0.75 | 0.125 | 1    | 0.125 |
| class 2 | 2.3  | 0.575 | 4.55 | 0.35  | 0.6  | 0.1   | 0.75 | 0.125 |
| class 3 | 0.7  | 0.175 | 0.5  | 0.25  | 3.3  | 0.55  | 3.45 | 0.575 |

# An example

# Confusion Matrices of Classification on Military Targets

|  | AAV | DW | HMV |
|-----|-----|-----|-----|
| AAV | 29 | 2 | 1 |
| DW | 0 | 18 | 8 |
| HMV | 0 | 2 | 23 |



Acoustic (75.47%, 81.78%)



Seismic (85.37%, 89.44%)



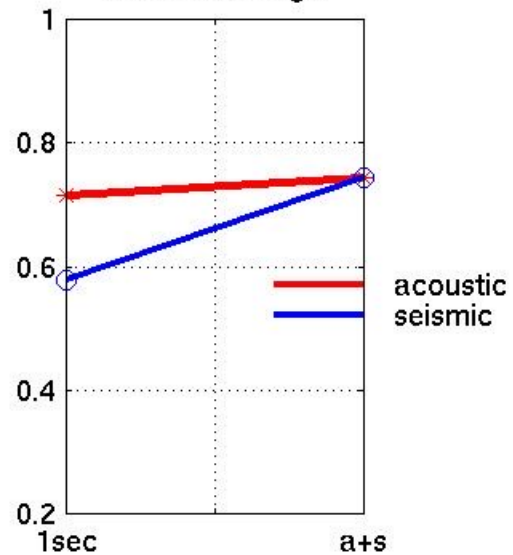Multi-modality fusion (84.34%)



Multi-sensor fusion (96.44%)

THE UNIVERSITY OF TENNESSEE KNOXVILLE

Acoustic

Seismic

Multi-modal

# Reference

- For details regarding majority voting and Naïve Bayes, see

http://www.cs.rit.edu/~nan2563/combining_classifiers_notes.pdf