

Combining Classifiers

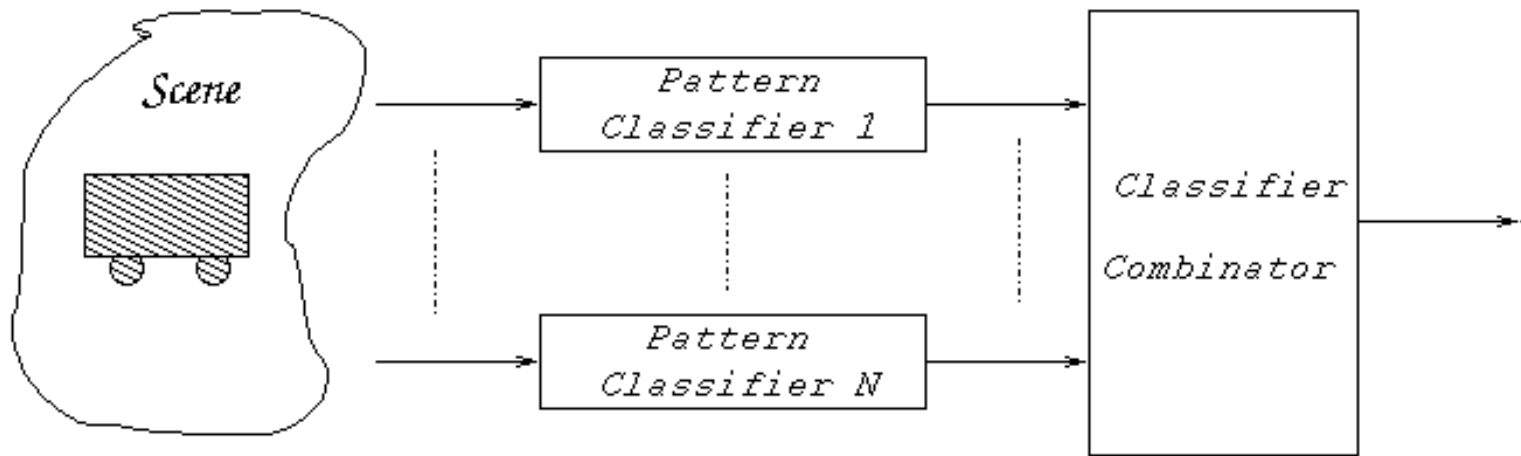
Sargur Srihari

srihari@cedar.buffalo.edu

Synonyms for the Topic

- “Mixture of Experts”
- “Ensemble Classifiers”
- “Modular Classifiers”
- “Pooled Classifiers”
- “Combining Models”

Typical Scenario



Application in Biometrics: Several Modalities

Two Approaches

1. Component Classifiers with Discriminant Functions

Use a principled statistical approach

2. Component Classifiers without Discriminant Functions

Several heuristics are available

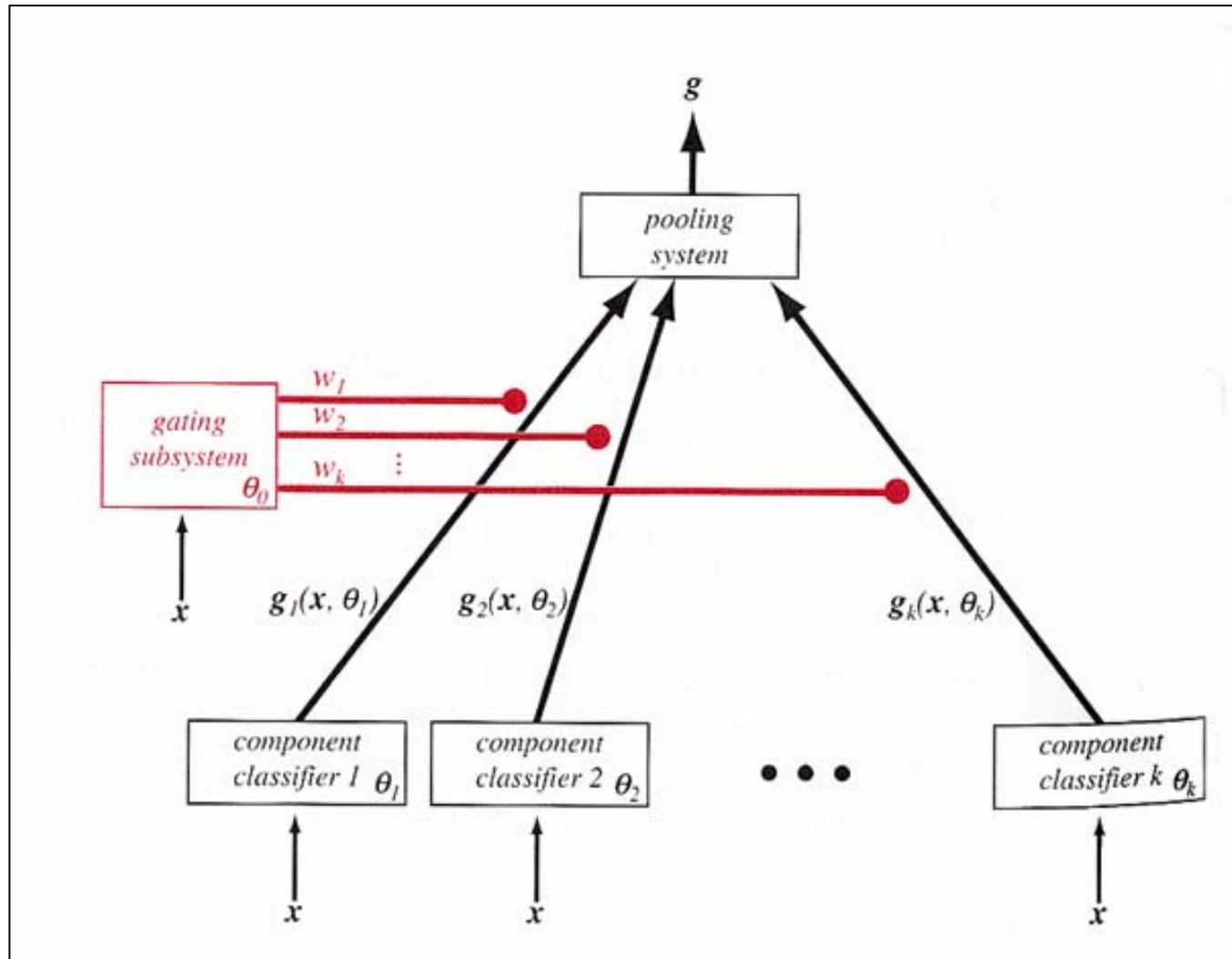
Mixture Model of Combination

- Each classifier tuned to its sub-distribution
- Each classifier outputs discriminant values
- Discriminant values are combined using weights for the components

Classifier Ensemble

- Each component classifier trained in a different region of feature space
- Each component classifier provides probability estimates

Architecture for Combining Classifiers



Mixture Model

- k classifiers corresponding to k component densities
- Classifier/component density (process indexed by r) chosen with probability

$$P(r | x, \theta_r^0)$$

- Where θ_r^0 is a parameter vector that describes the state of the process

Probability of producing a category label y

Mixture distribution

$$P(y|x, \Theta^0) = \sum_{r=1}^k P(r|x, \theta_r^0) P(y|x, \theta_r^0)$$

$$\Theta^0 = [\theta_0^0, \theta_1^0, \dots, \theta_k^0]$$

Vector of all relevant parameters

Discriminant Functions of Component Classifier r

$$\sum_{j=1}^c g_{rj} = 1 \quad \text{for all } r.$$

Conditional Mean of $P(r | x, \theta_r^0)$

Gating subsystem with weight parameters

(All discriminant values of component classifier r are multiplied by a scalar w_r) :

$$\mu = \mathcal{E}[y | x, \Theta] = \sum_{r=1}^k w_r \mu_r$$

Parameter Estimation

Find parameters that maximize the log-likelihood function for n training patterns x^1, \dots, x^n in \mathcal{D} :

$$l(\mathcal{D}, \Theta) = \sum_{i=1}^n \ln \left(\sum_{r=1}^k P(r|x^i, \theta_0) P(y^i|x^i, \theta_r) \right)$$

Use Gradient Descent on the Parameters

$$\frac{\partial l(\mathcal{D}, \Theta)}{\partial \mu_r} = \sum_{i=1}^n P(r|y^i, x^i) \frac{\partial}{\partial \mu_r} \ln[P(y^i|x^i, \theta_r)]$$

for $r = 1, \dots, k$

Gradient Descent moves prior probabilities to posterior probabilities

$$\frac{\partial l(\mathcal{D}, \Theta)}{\partial \mu_r} = \sum_{i=1}^n P(r | y^i, x^i) - w_r^i$$

Winner Take All

- Use the decision of the single component classifier that is most confident
- Has the largest discriminant value,
- Sub-optimal, easy to implement
- Can work well if classifiers are experts in separate regions of the input space

g_{rj}

Component Classifiers Without Discriminants

- Wish to create ensemble classifier from highly trained component classifiers some of which may not compute discriminant functions
- Output values
 - Analog values, e.g., neural network
 - Rank order, e.g., k-nearest neighbor
 - One of c , e.g., rule-based system

Conversion of Analog values \tilde{g}_i output by component classifiers
to Discriminant Values using SOFTMAX transformation

$$g_i = \frac{e^{\tilde{g}_i}}{\sum_{j=1}^c e^{\tilde{g}_j}}$$

Conversion of Rank Order values to Discriminant Values

Assume discriminant function is
linearly proportional to rank order on the list

Resulting g_i should be normalized
so that they sum to one

One of c :

If the output is a one-of- c representation, in which a single
Category is identified, we let $g_j=1.0$ and 0.0 otherwise

Conversion to Discriminant Values

The table gives a simple illustration of these heuristics.

Analog value		Rank order		One-of- c	
\tilde{g}_i	g_i	\tilde{g}_i	g_i	\tilde{g}_i	g_i
0.4	0.158	3rd	$4/21 = 0.194$	0	0.0
0.6	0.193	6th	$1/21 = 0.048$	1	1.0
0.9	0.260	5th	$2/21 = 0.095$	0	0.0
0.3	0.143	1st	$6/21 = 0.286$	0	0.0
0.2	0.129	2nd	$5/21 = 0.238$	0	0.0
0.1	0.111	4th	$3/21 = 0.143$	0	0.0