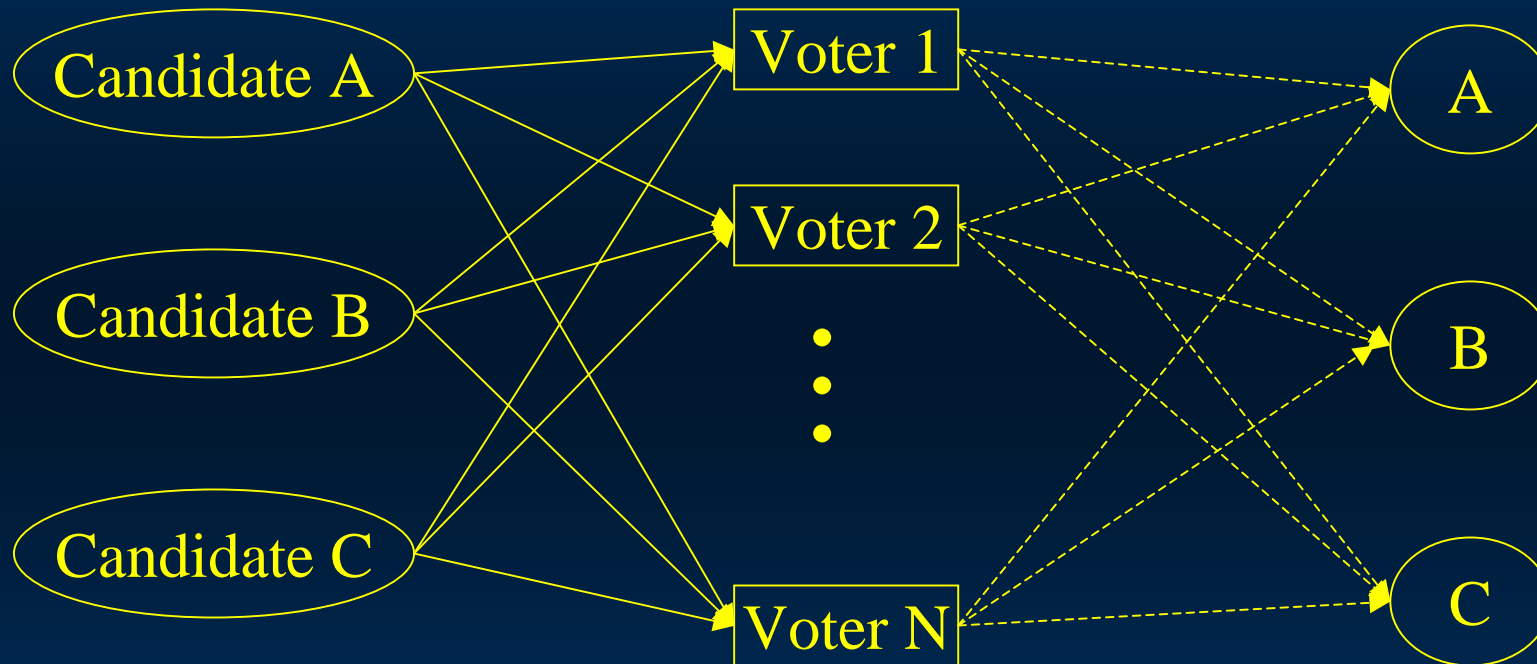


Classifier Combination for Biometrics

Voting Example (1)



- Each candidate is a class.
- Each voter is a classifier.
- The voting problem is to define the type of voter's output and to combine the results of the vote.

Voting Example (2)

Majority voting:

each voter can vote for one candidate; candidate with biggest number of votes wins.

Suppose:

200 voters – $A > C > B$

195 voters – $B > C > A$

20 voters – $C > B > A$

Is C the worst candidate?

Consider candidates in pairs:

C and A : 215 voters for C and 200 for A

C and B : 220 voters for C and 195 for B

Voting Example (3)

Borda (18th century) count:

Each voter gives to the first choice candidate 2 points, 1 point to second choice and 0 for third choice candidate.

In our example:

A – 400 points, B – 410 and C – 435 points.

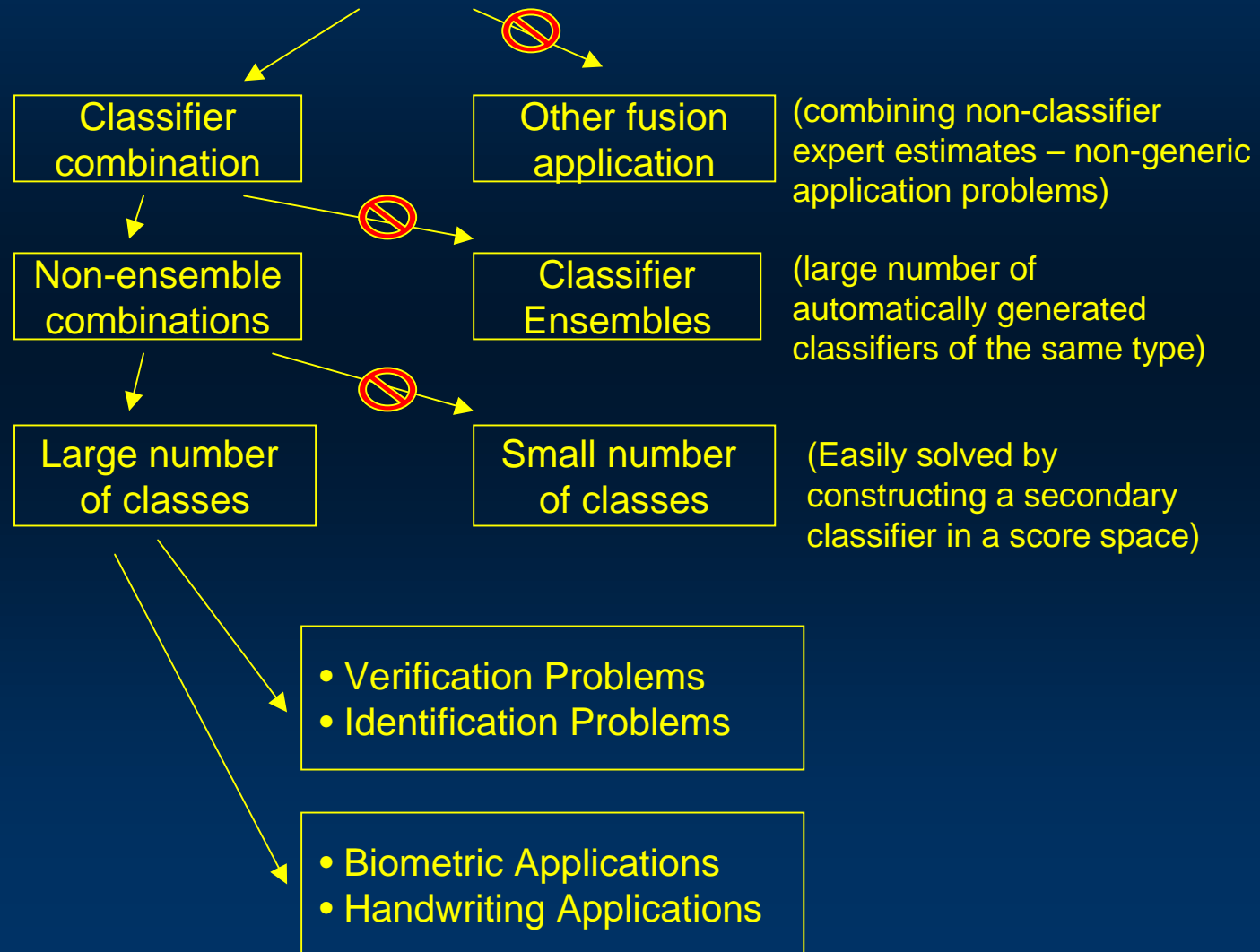
Approval voting :

Each voter gives one point to every candidate he approves, and no points for others.

Voting Example (4)

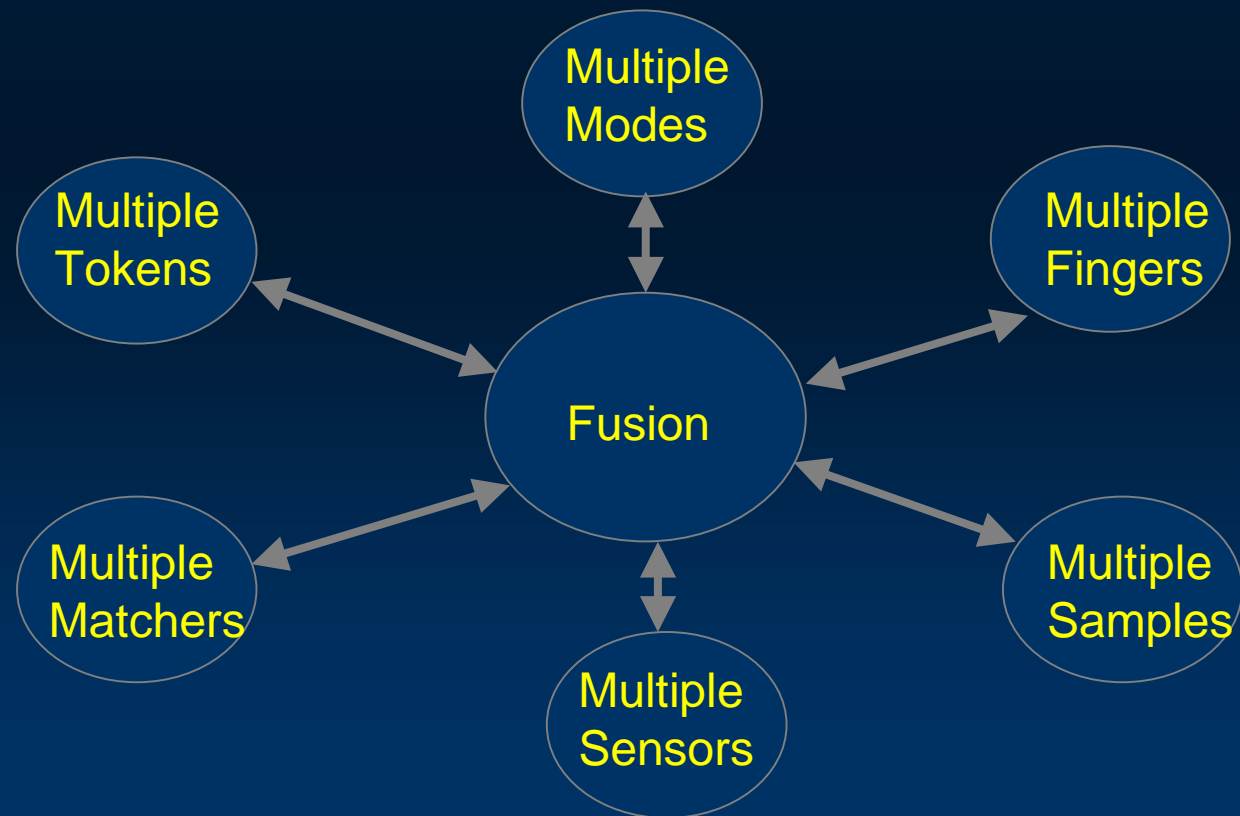
- Additional modifications:
 - See http://en.wikipedia.org/wiki/Category:Voting_systems
 - Multi-stage votes
 - Confidence vote
- Main problem with defining voting methods:
What is the 'true class' or 'true candidate' ?

Combination Problems



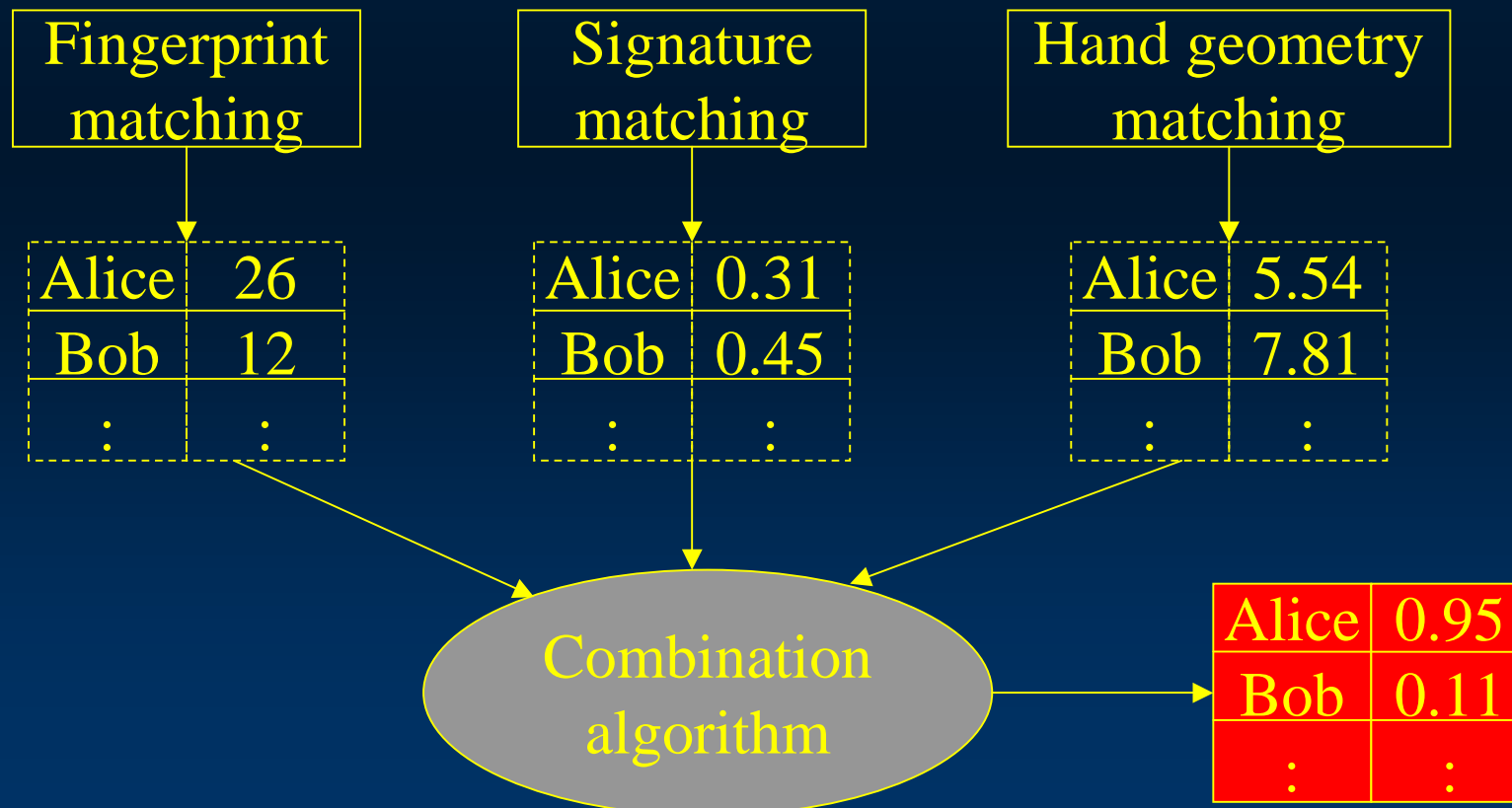
Multi-Biometrics

- **Single biometric may not be**
 - Secure
 - Provide sufficient coverage
 - Accurate
 - Efficient



Combination of biometric matchers

- Criteria for classifier combination- performance on training/testing sets.
- Classifiers can produce different types of scores – integers, floats



Types of classifiers and type conversion

	Type I	Type II	Type III
Output	One class is chosen	Ranking	Confidence scores
Example of scores	$S(C_1)=1$ $S(C_2)=0$ $S(C_k)=0$ $S(C_n)=0$	$S(C_1)=n-1$ $S(C_2)=n-2$ $S(C_k)=n-k$ $S(C_n)=0$	$S(C_1)=.4$ $S(C_2)=.3$ $S(C_k)=.05$ $S(C_n)=0$

- Combination methods usually accept input of one type.
- If combined classifiers produce different type of input, then conversion is required.

Combination algorithms

	Output values	Combination algorithms
Type I	The single class	Majority vote
Type II	Class ranking	Borda count
Type III	Class ranking with confidence score	Sum rule, Neural network.

Combination for identification task

- Typical combination method – sum rule:

$$S(w_i) = \sum_j \alpha_j s_j(w_i)$$

- The person is identified as c , where

$$c = \arg \max_i S(w_i)$$

Combination for verification task

- It is assumed that there are two classes: person to be verify and all other persons.
- Typical combination rule for verification: find

$$S(w) = \sum_j \alpha_j s_j(w)$$

and see if

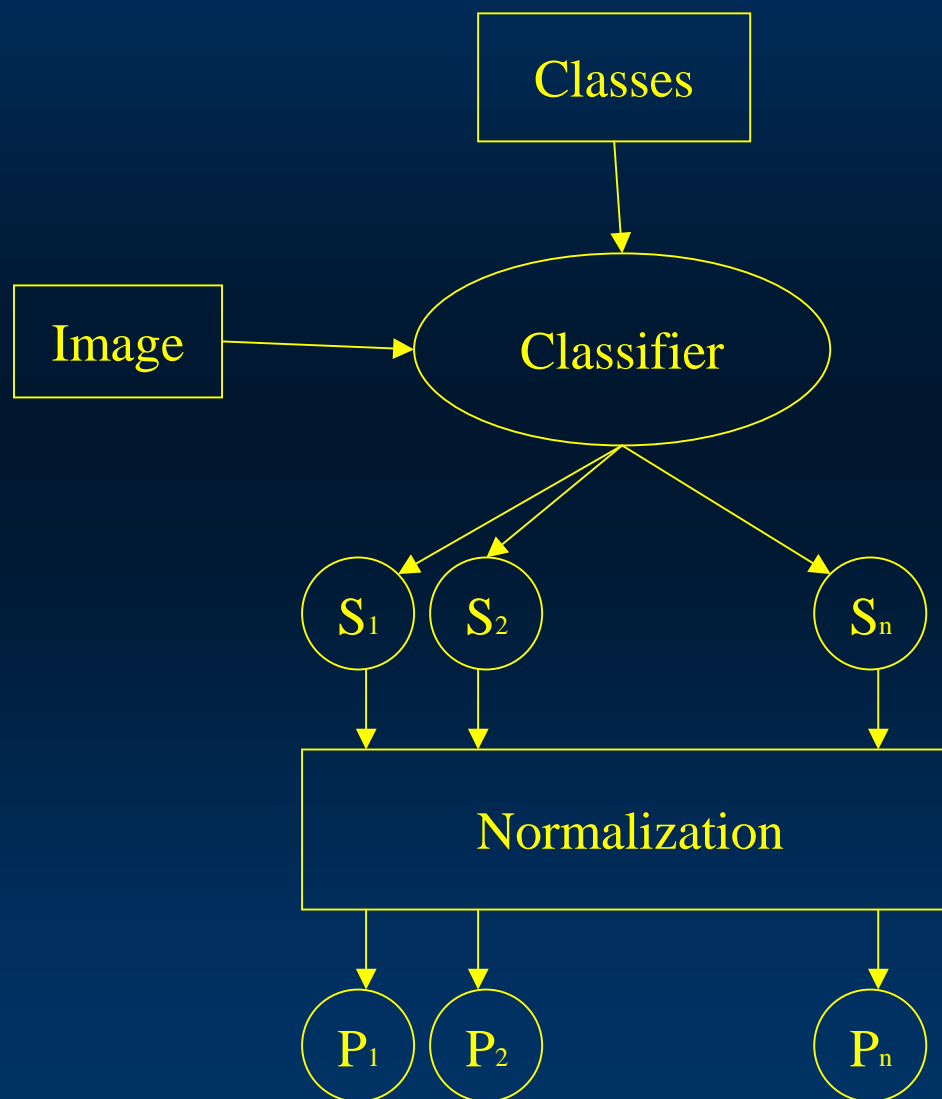
$$S(w) > \theta$$

for some predefined threshold θ

Classifiers' score conversion

- Scores reflect
 - Distance measures, confidence values, and beliefs
 - Typically, represent the distance of a pattern from its prototype in *some* feature space using *some* metric
- **Scores across classifiers are at a different scale**
- **Integration of classifier output with other modules becomes ad hoc**
- **Need for the probability of correctness**

Score normalization:



Class score example (1)

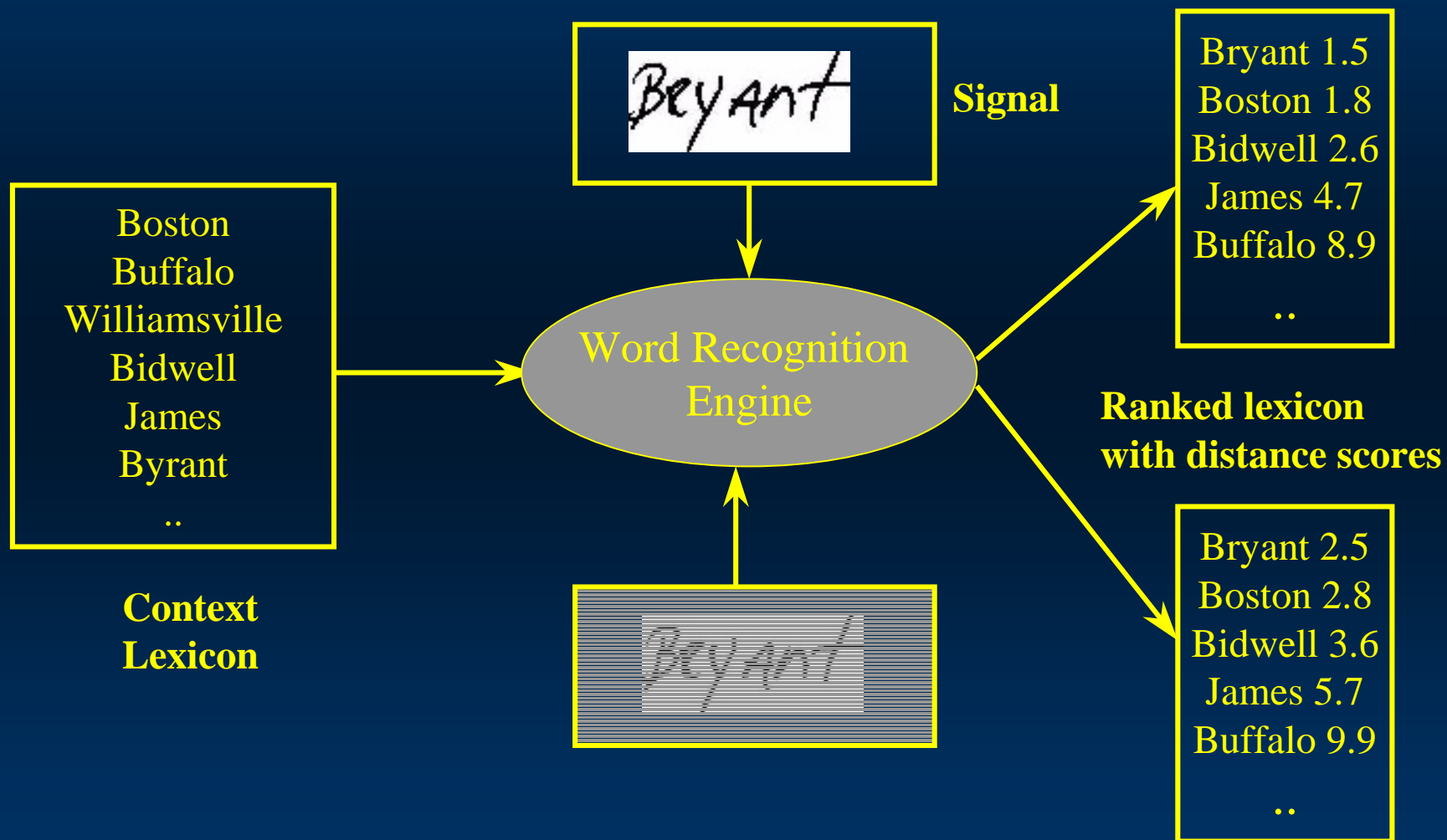
- Teacher A and Teacher B want to evaluate the proficiency of n students in Math.
- Both teachers give an exam to the students and evaluate their responses.
- Student S scores the highest on A's evaluation: 80%
- Student S scores the highest on B's evaluation: 90%
- In the analogy:
 - *Teachers are classifiers*
 - *Students are the pattern classes*
 - *Input patterns are the exams*

Class score example (2)

- *Question 1:* Is S the most proficient in Math in the class?
- *Question 2:* Is the opinion of B about Omega's proficiency stronger than the opinion A has of Omega's proficiency?
- *Question 3:* Given that A gave 80% as the highest score, what is the probability that he is correct in choice of best student?

- *Answer 1:* If S can consistently rank first over many exams of "same" difficulty
- *Answer 2:* Study the grading behavior of the teachers over many exams
- *Answer 3:*
 - ✓ Use rankings, not scores.
 - ✓ Derive the probability of correctness from score

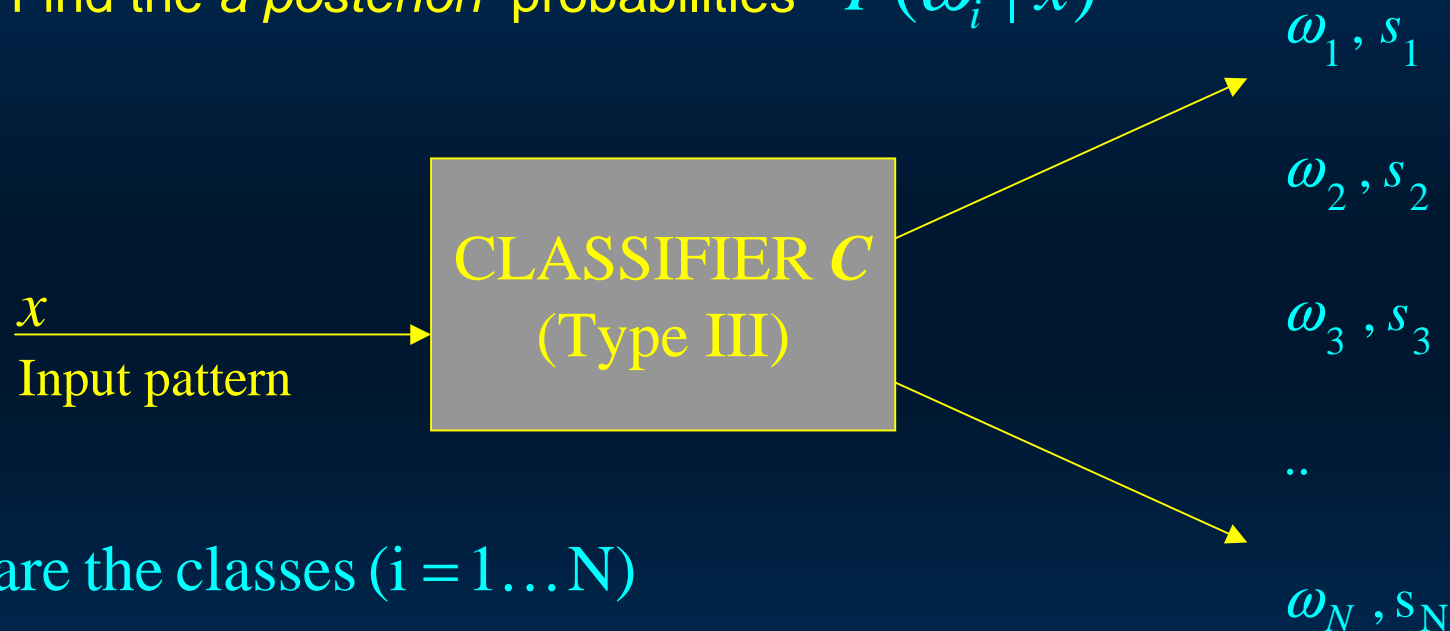
Using rankings instead of scores



Rankings can be more reliable than scores.

Deriving probabilities (1)

- Find the *a posteriori* probabilities $P(\omega_i | x)$



ω_i are the classes ($i = 1 \dots N$)

$s_i \in [0,1]$ are the scores returned by the classifier

$$\sum_{i=1 \dots N} P(\omega_i | x) = 1$$

- Instead of $P(\omega_i | x)$ consider $P(\omega_i | s_1, \dots, s_c)$ or even $P(\omega_i | s_i)$

Deriving probabilities (2)

- Two probability density functions could be used for score normalization:

$$p(C(x, w_i) = s | t(x) = w_i)$$

- score distribution when truth of the input pattern is the same as class w_i

$$p(C(x, w_i) = s | t(x) \neq w_i)$$

- score distribution when truth of the input pattern is different from the class w_i

Deriving probabilities (3)

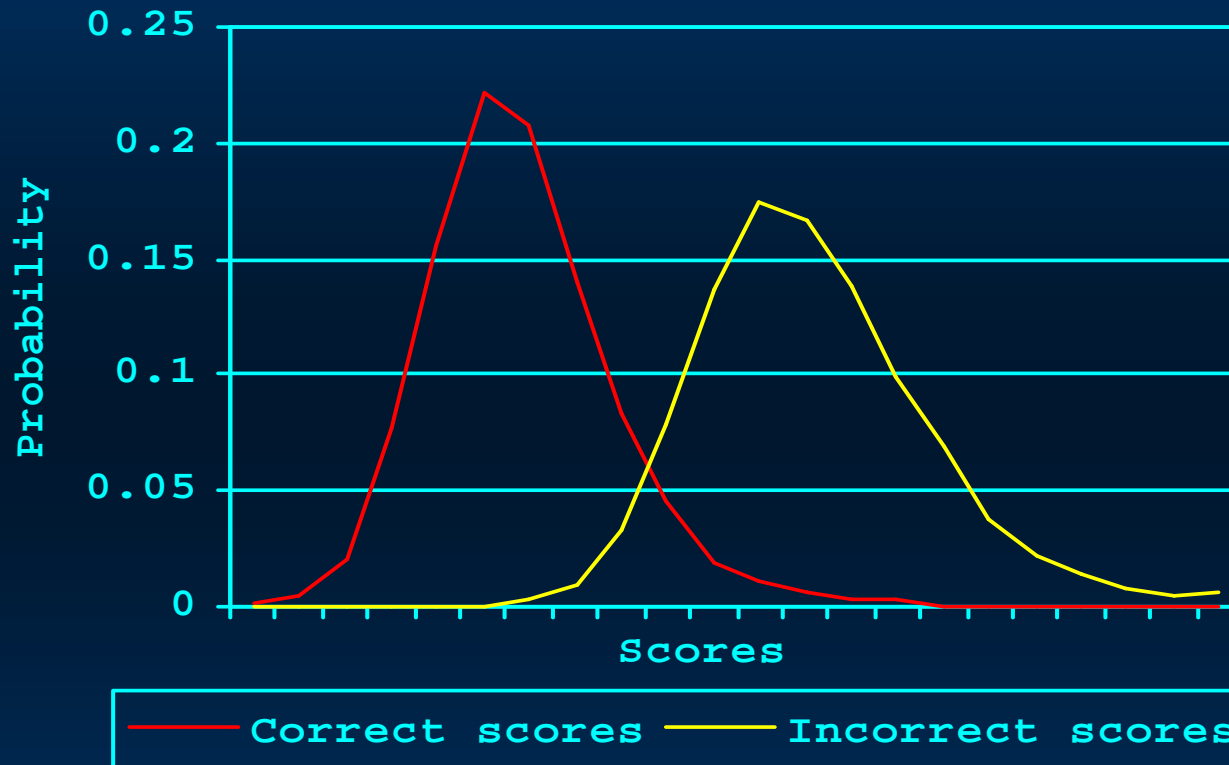
Using these two distributions we can map the recognizer's score into probability that word is recognized correctly:

$$P(t(x) = w_i | C(x, w_i)) = \frac{p(C(x, w_i) | t(x) = w_i)P(t(x) = w_i)}{p(C(x, w_i))}$$

and

$$p(C(x, w_i)) = p(C(x, w_i) | t(x) = w_i)P(t(x) = w_i) + p(C(x, w_i) | t(x) \neq w_i)P(t(x) \neq w_i)$$

Histograms to derive probabilities



Correct scores distribution is same as $p(C(x, w_i) = s | t(x) = w_i)$

Incorrect scores distribution is same as $p(C(x, w_i) = s | t(x) \neq w_i)$

Using logistic function for score normalization

Logistic function:
$$f(s) = \frac{1}{1 + e^{as+b}}$$

- $f(s)$ has a range between 0 and 1. Thus it is convenient candidate for modeling probabilities.
- We need to find parameters a and b , such that $f(s)$ most closely coincides with $P(\omega | s)$
- Finding parameters statistically from given training samples is called *regression*.

Summary of Approaches to Score Normalization

1. Use rankings instead of scores.
2. “Linearly” output scores to some specified range.
3. Histogram method.
4. Model probability of correct match by regression methods.
5. Rely on higher level functions for normalization, e.g., neural networks.

Combination methods (1)

- Majority vote : count the number of votes given for all classes and choose class with biggest number of votes.
 - little information is used, might want to use ranking information also (as in election example)
 - if number of classifiers is small, there will be a lot of ties
- Borda count : sum up the rankings for each class, choose class with best sum.
 - frequently used
 - sometimes wan to use score information
- Sum rule: sum up scores for each class, choose class with best sum.
 - requires score normalization
 - might be worse than Borda count

Combination methods (2)

- Logistic function : model the probability that class is correct given output scores of classifiers:

$$P(t(x) = w_i | s_i^{(1)}, \dots, s_i^{(M)}) \approx$$
$$\approx f(s_i^{(1)}, \dots, s_i^{(M)}) = \frac{1}{1 + e^{a_1 s_i^{(1)} + \dots + a_M s_i^{(M)} + b}}$$

- Neural network : same as for logistic function – model probability of correct classification:

$$P(t(x) = w_i | s_i^{(1)}, \dots, s_i^{(M)})$$

Note, that neural network can model arbitrary function.

Summary of Combination Methods

Approaches	Logic based	“Try them all”	Statistical
Description	-Assumptions on the meaning of combined data -The combination algorithm is a predetermined rule	Try few predetermined rules; choose one with best performance	The combination function is derived using training data and machine learning algorithms
Ease of use	Average	Easy	Difficult
Training data requirements	Low	Average	High
Optimality of combination	No	Somewhat	Yes

Example of logic based approach – Dempster-Shafer Theory

$P(X)$ - the power set of X

$m : P(X) \rightarrow [0,1]$ - basic belief assignment

Belief: $bel(A) = \sum_{B|B \subseteq A} m(B)$ Plausibility: $pl(A) = \sum_{B|B \cap A \neq \emptyset} m(B)$

Dempster's combination rule:

$$m_{1,2}(A) = (m_1 \oplus m_2)(A) = \frac{1}{1-K} \sum_{B \cap C = A} m_1(B)m_2(C)$$

$$K = \sum_{B \cap C = \emptyset} m_1(B)m_2(C)$$

Not optimal:

- basic belief assignments are heuristically chosen
- many other similar rules were proposed claiming superior performance
- assumes statistical independence of combined events

Example of “Try them all” approaches

Kittler et al., “On Combining Classifiers” , 1998:

- 6 rules are justified under different assumptions:

-Sum rule $S_i = f(s_i^1, \dots, s_i^n) = s_i^1 + \dots + s_i^n$

-Product rule $S_i = f(s_i^1, \dots, s_i^n) = s_i^1 \times \dots \times s_i^n$

-Max rule $S_i = f(s_i^1, \dots, s_i^n) = \max(s_i^1, \dots, s_i^n)$

-Min rule $S_i = f(s_i^1, \dots, s_i^n) = \min(s_i^1, \dots, s_i^n)$

-Median rule $S_i = f(s_i^1, \dots, s_i^n) = \text{median}(s_i^1, \dots, s_i^n)$

-Majority vote

$$S_i = f(s_i^1, \dots, s_i^n) = \sum_j v_i^j, \quad v_i^j = \begin{cases} 1, & \text{if } s_i^j > s_i^k \quad \forall k \\ 0, & \text{otherwise} \end{cases}$$

s_i^j - score assigned to class i
by the classifier j

Somewhat optimal:

- choose best performing rule
- no confidence that chosen rule is close to optimal
- multiple published results show that different rules can be best in different problems

Our Method

(Work done with Ultra-Scan)

- Sort the hypercubes

$$\mathbf{V}_S = \{V_i : P(M | \mathbf{x} \in V_{i+1}) > P(M | \mathbf{x} \in V_i), i = 1..N_V\}$$

- Given a desired FAR, p_{fm} , find N_{fm} , such that

$$p_{fm} = \sum_{i=1}^{N_{fm}} P(\mathbf{x} | NM)$$

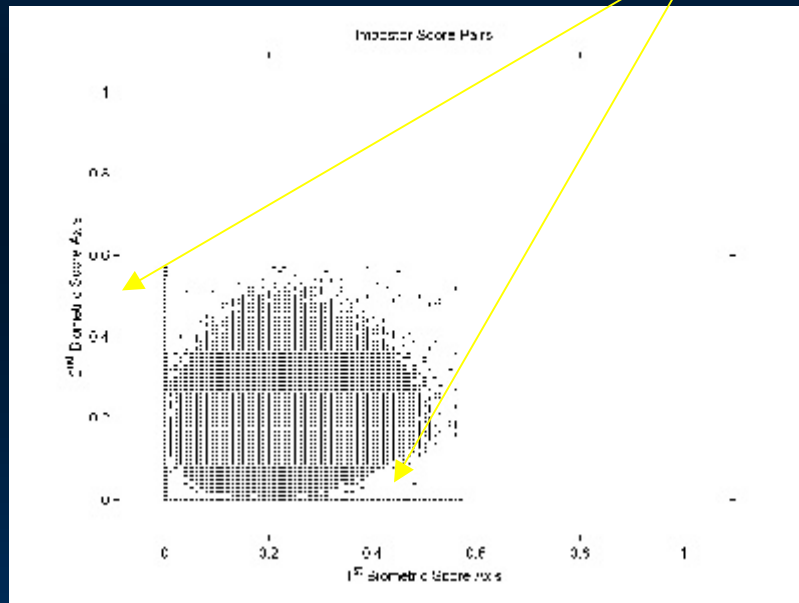
- Optimal Decision regions for desired FAR

$$R_M = \bigcup_{i=1}^{N_{fm}} V_i \quad \text{and} \quad R_{NM} = \bigcup_{i=N_{fm}+1}^{N_V} V_i$$

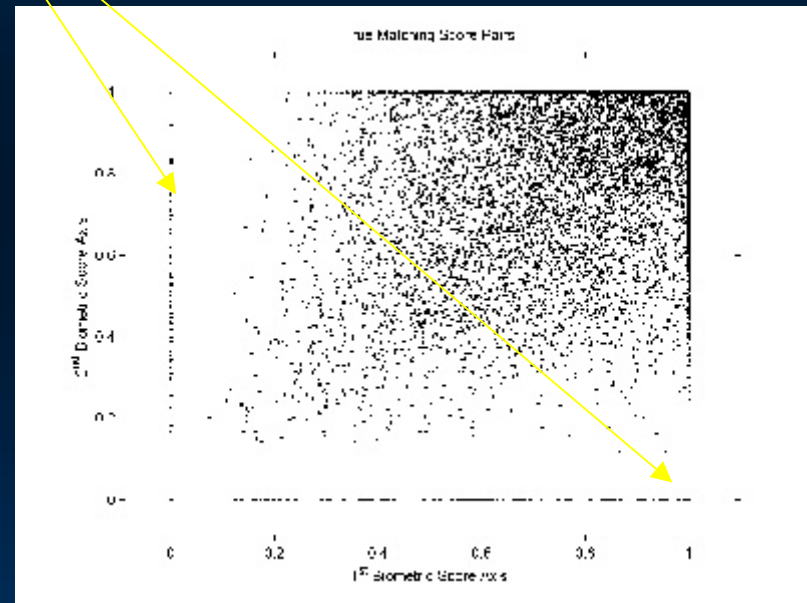
Two Fingerprint Fusion Example

Matching Score Pairs are Strongly Correlated

multimodal

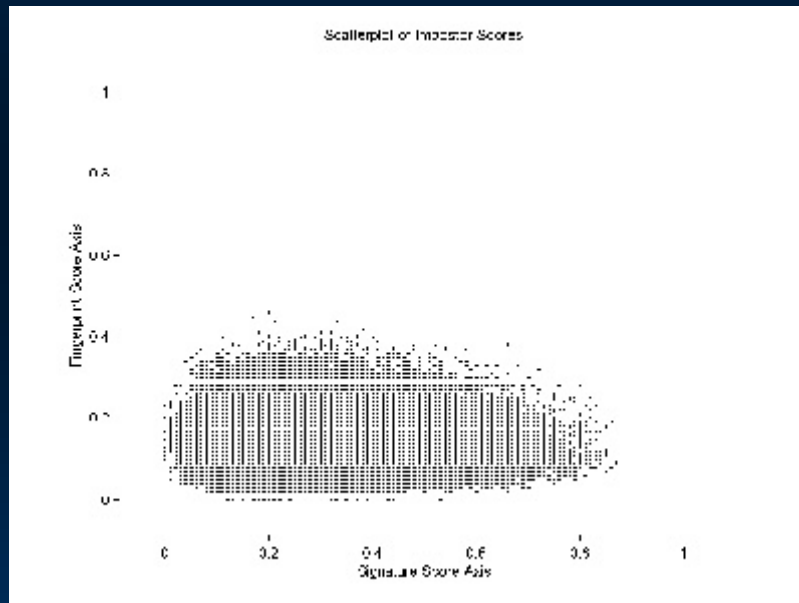


Impostor Score Pairs

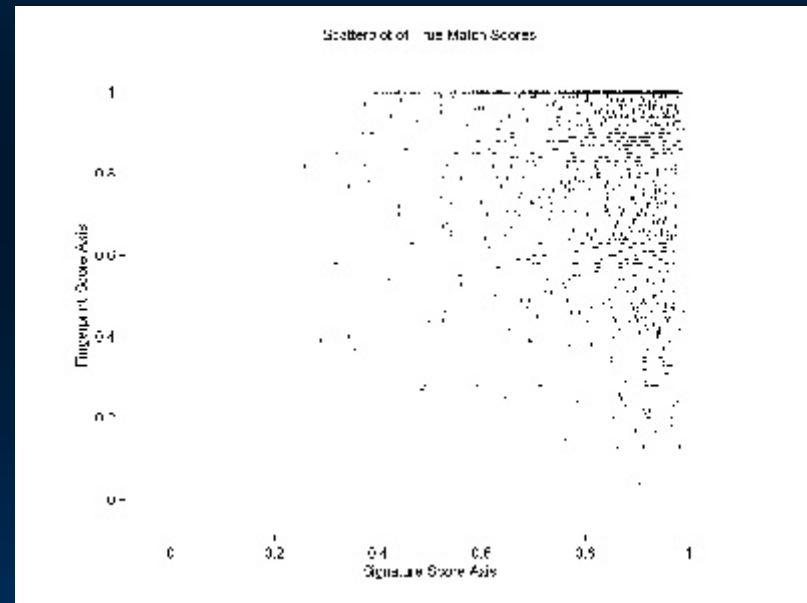


True Match Score Pairs

Signature & Fingerprint Fusion (Uncorrelated)



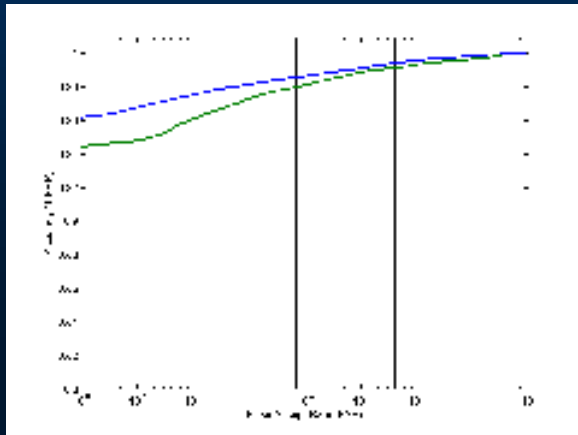
Impostor Score Pairs



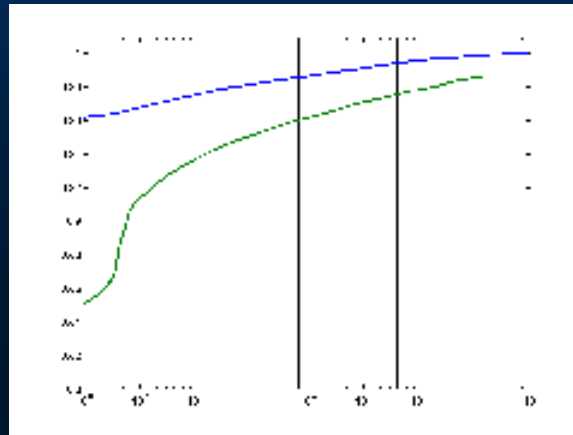
True Match Score Pairs

Fusion Methods

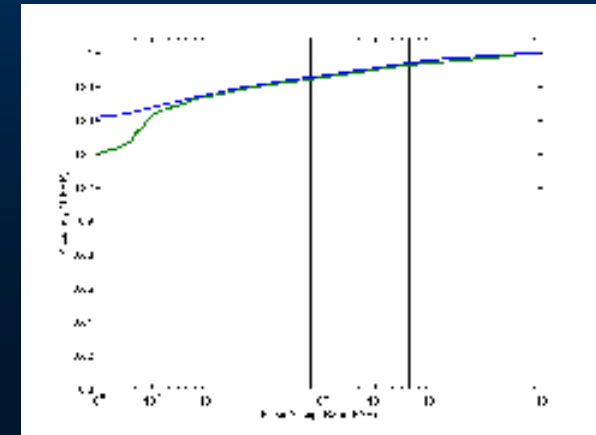
2 Fingers



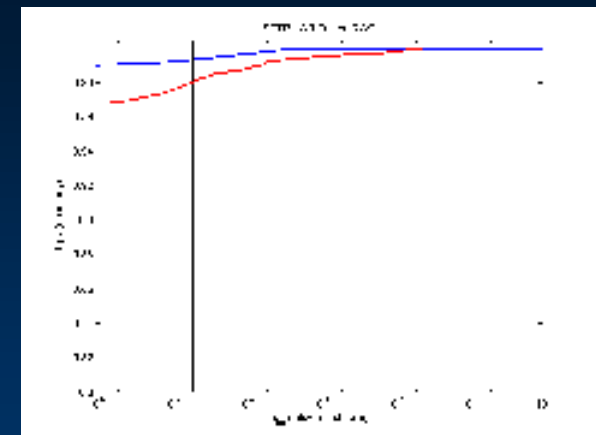
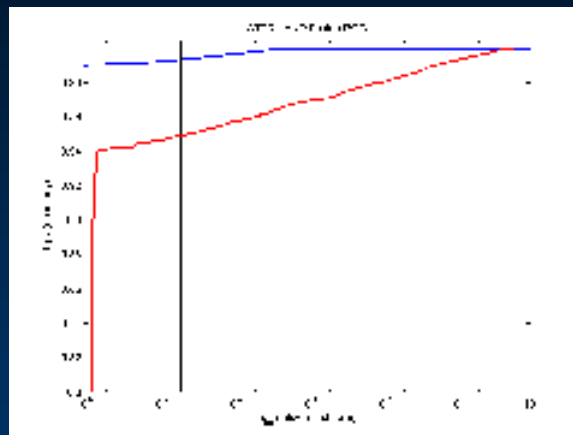
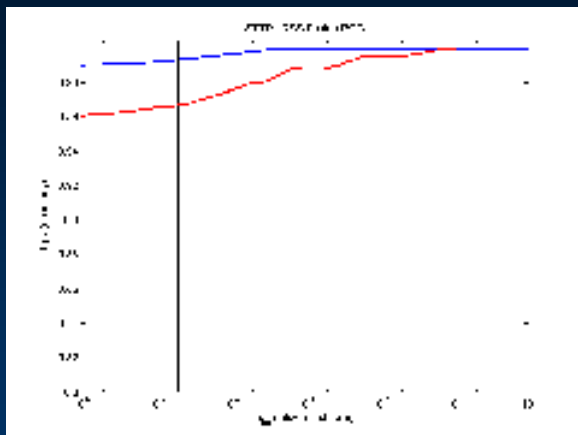
OR



AND

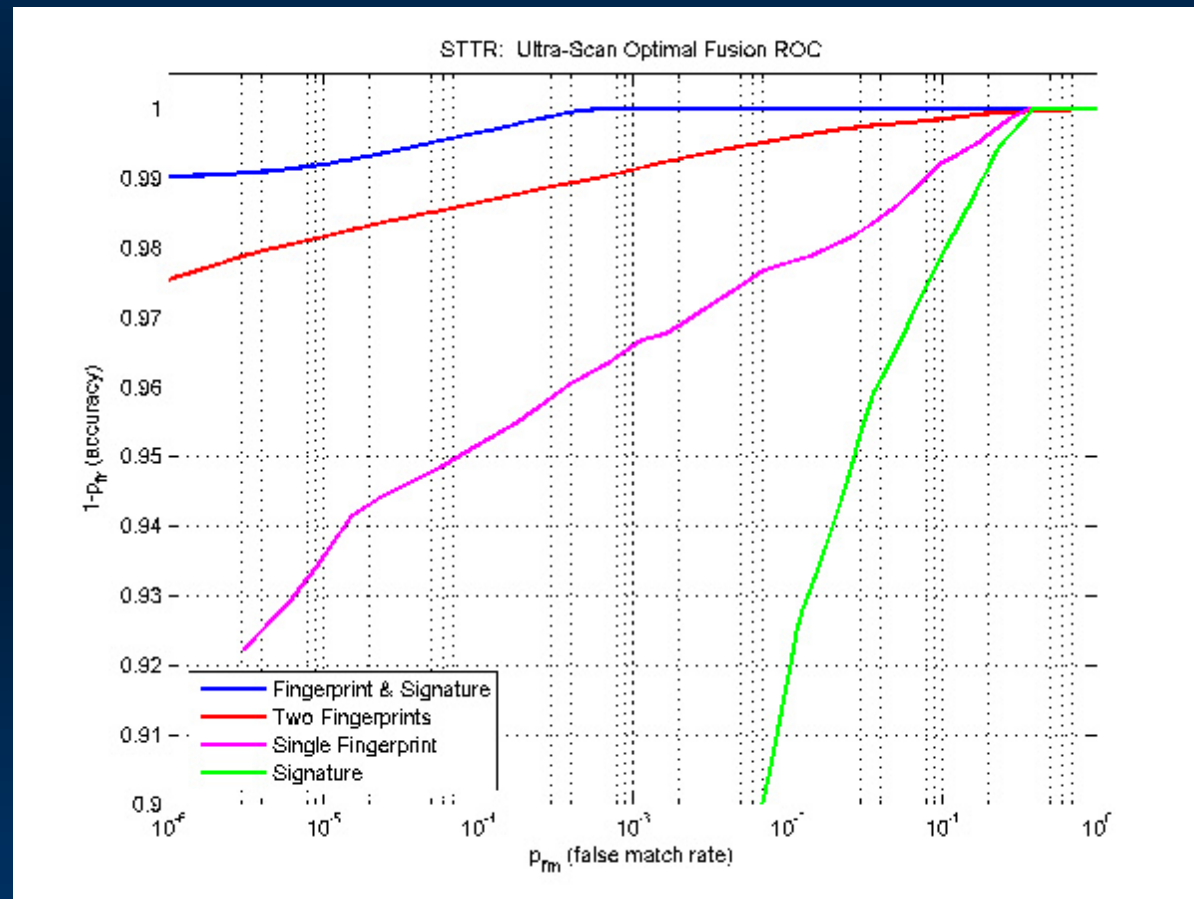


RMS



Fingerprint + Signature

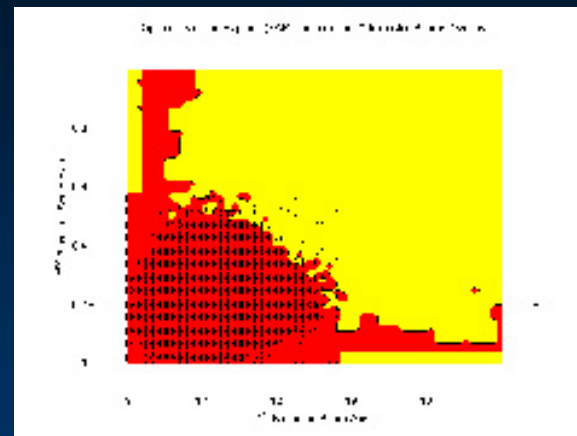
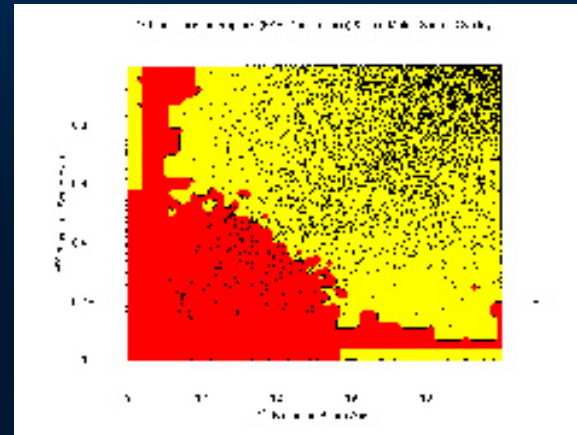
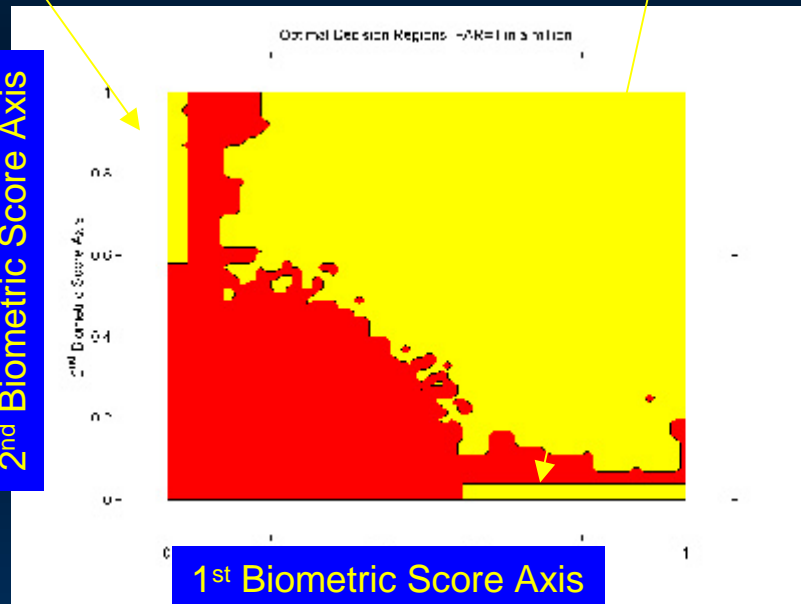
Optimal Fusion ROC Comparisons



2 Finger Fusion Algorithm 97.55% Accuracy @ Specified FAR of 1 in a Million

Note adaptation to multimodal match score zones

2nd Biometric Score Axis



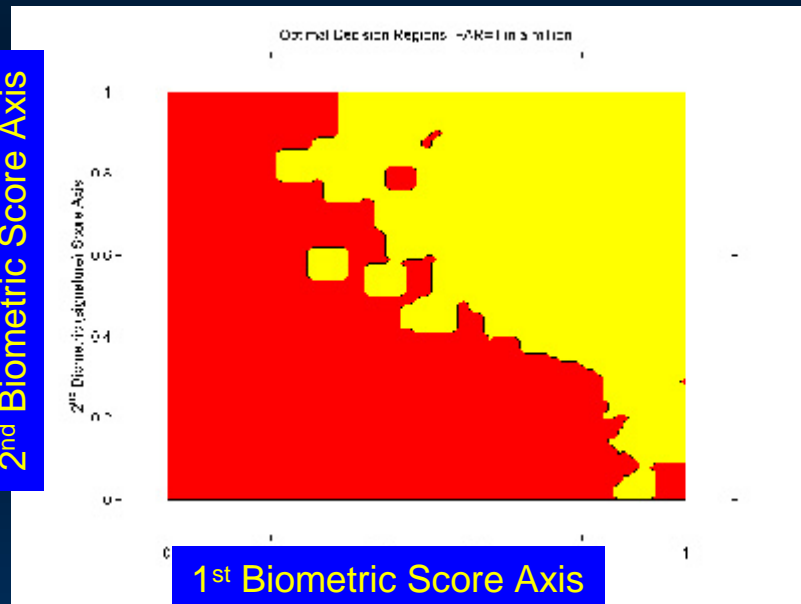
Match Zone

No-Match Zone

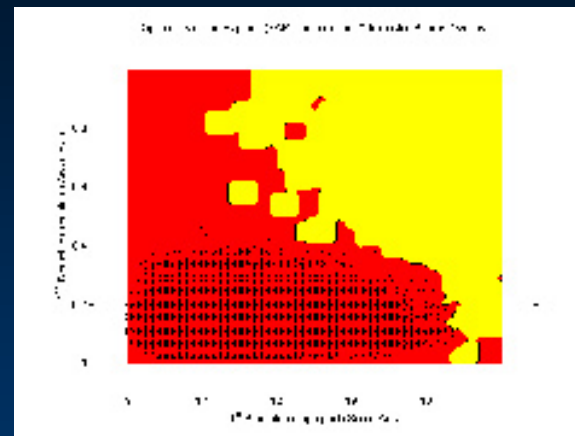
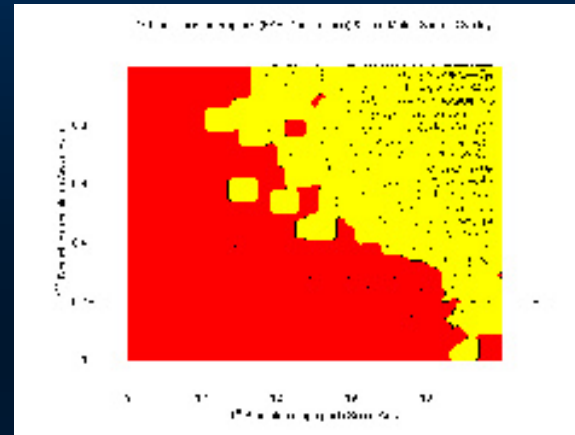
Fingerprint + Signature Fusion

99.04% Accuracy @ Specified FAR of 1 in a Million

2nd Biometric Score Axis



1st Biometric Score Axis

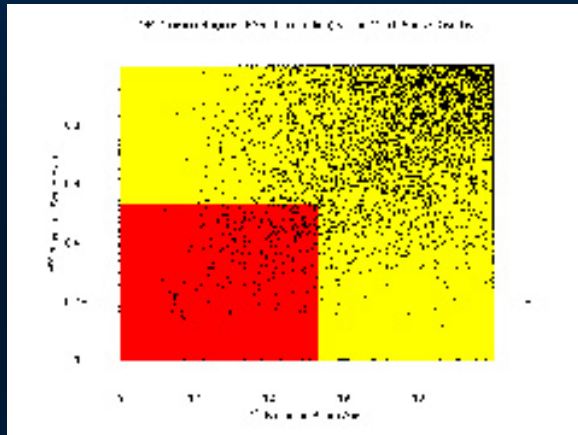


■ Match Zone

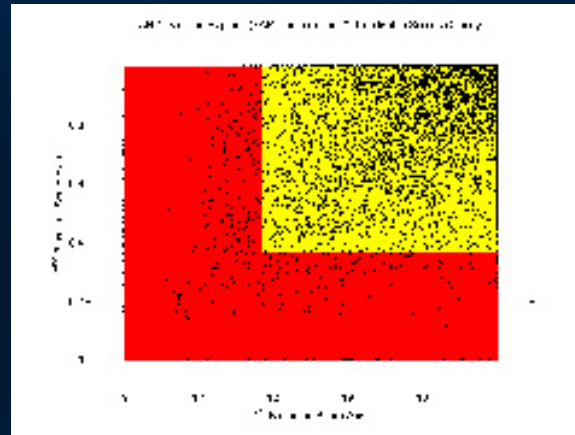
■ No-Match Zone

2 Finger Fusion at FAR of 1 in Million

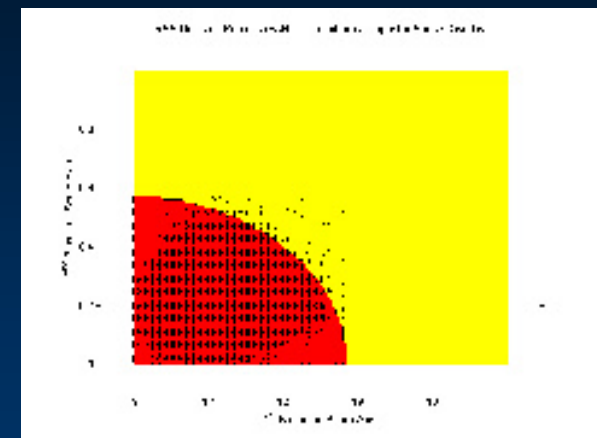
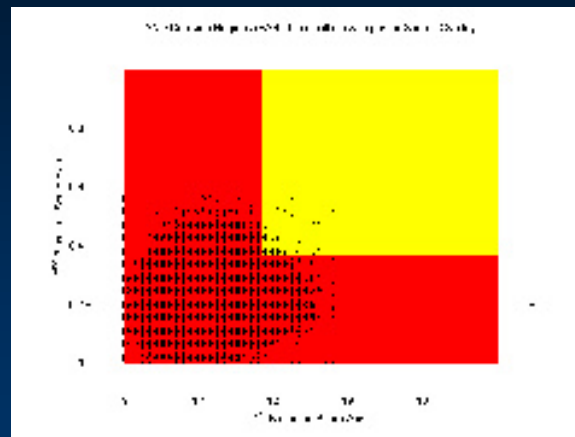
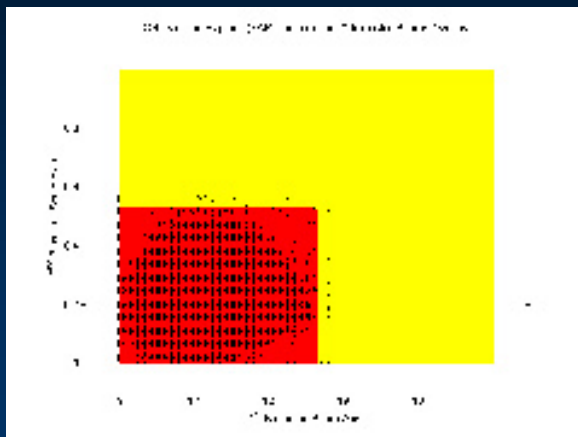
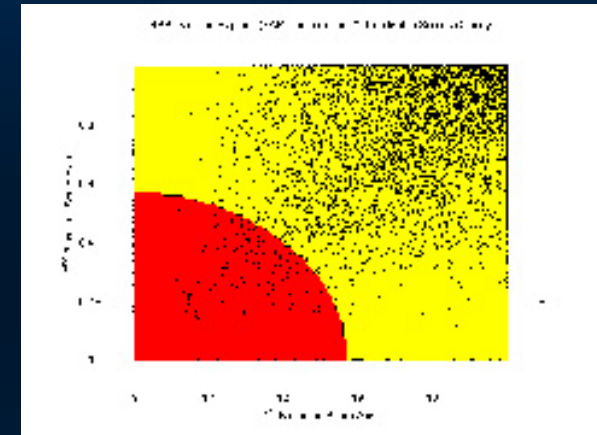
OR
96.11% Accuracy



AND
93.74% Accuracy



RSS
97.48% Accuracy



 Match Zone

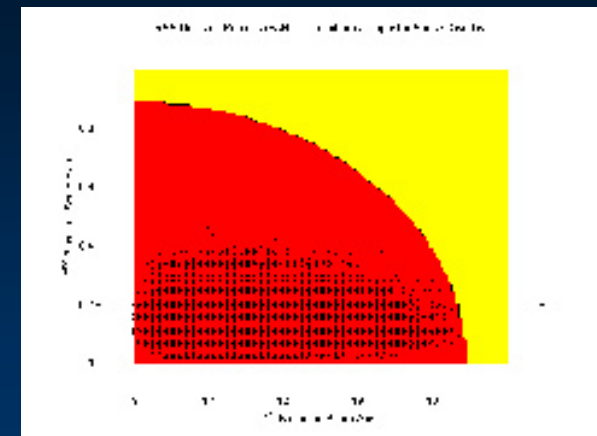
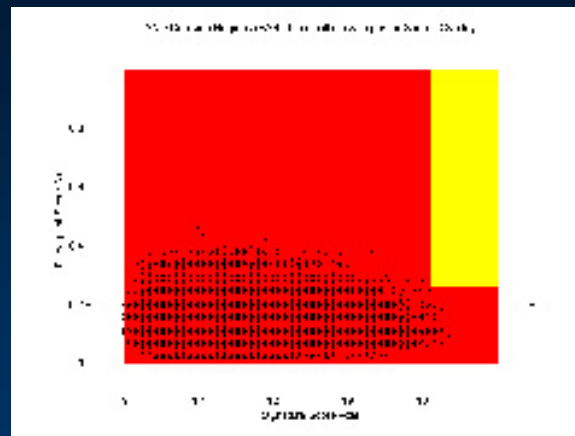
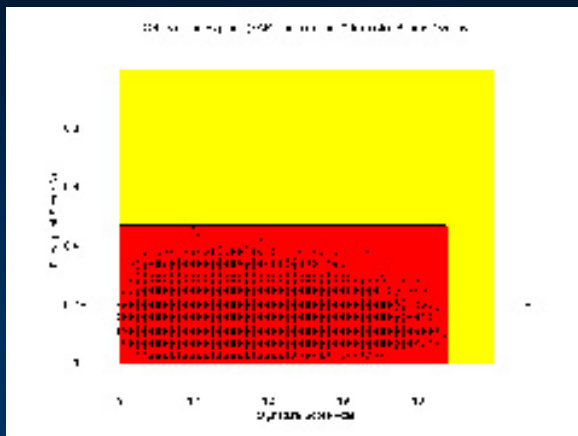
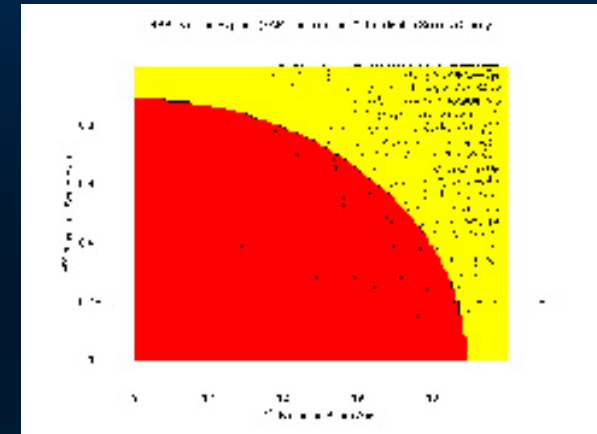
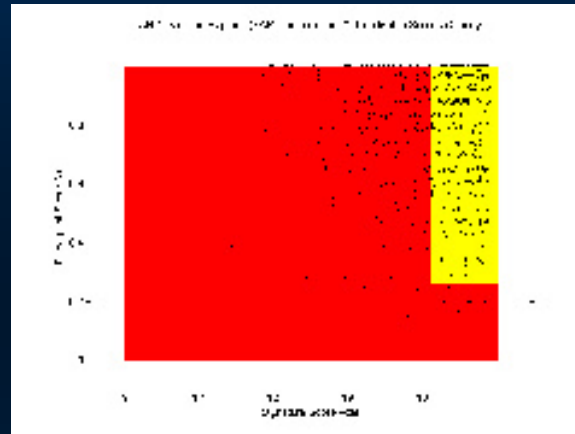
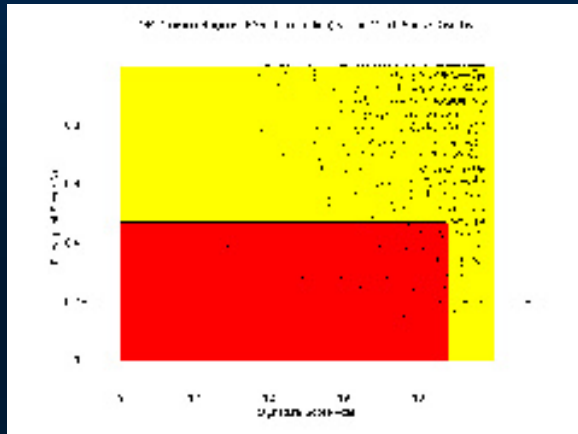
 No-Match Zone

Fingerprint + Signature Fusion at FAR of 1 in Million

OR: 96.85% Accuracy

AND: 62.91% Accuracy

RMS: 96.11% Accuracy



 Match Zone

 No-Match Zone

Statistical Approaches

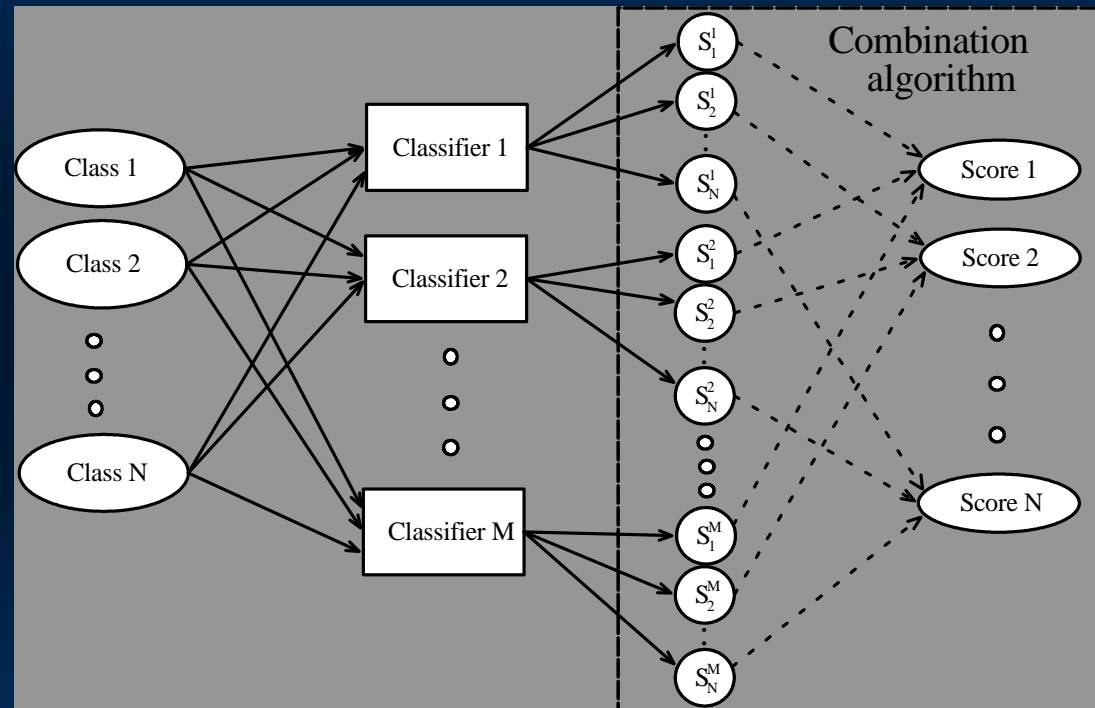
- Combination problem - a problem of learning combination algorithm from training samples
- A set of learning algorithms is chosen with unknown parameters
- The best parameters are found with respect to the cost function and training data
- It is possible to give an estimate on the proximity of found solution to the optimal

Advantages:

Universal approximation property guarantees the closeness to the optimal solution

But:

Need to properly choose cost function and avoid overfitting



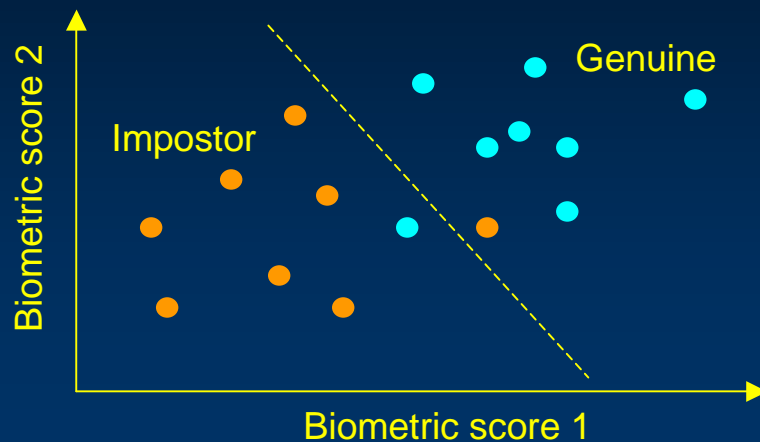
Bayesian Risk Minimization Classifier

- Treat a set of scores as a feature vector, and solve a pattern classification problem with two classes: genuine and impostor.
- Neyman Pearson Likelihood Ratio

$$p(s_1, s_2 | imposter)$$

$$p(s_1, s_2 | genuine)$$

$$\frac{p(s_1, s_2 | genuine)}{p(s_1, s_2 | imposter)}$$

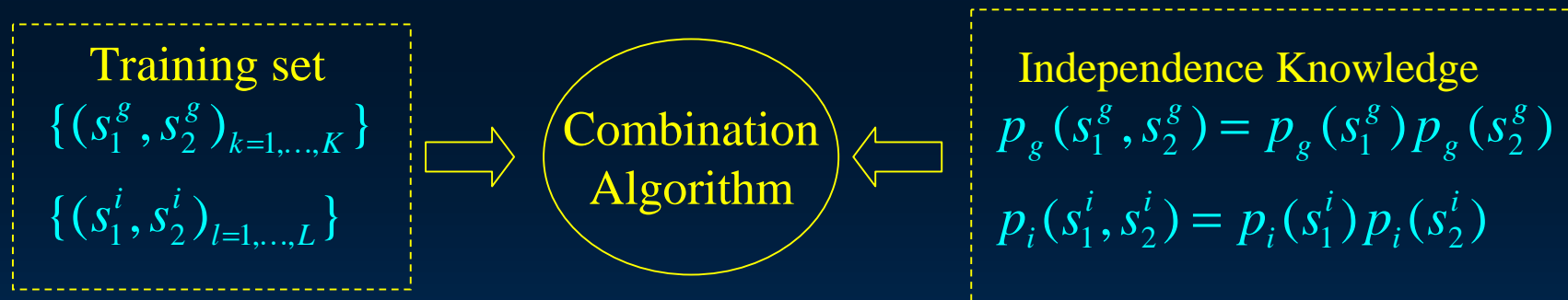


If pdfs are estimated accurately

Likelihood Ratio gives Optimal decision boundary

Using Independence Knowledge

- Large number of training samples
- Few biometric matchers
- Use Independence of matcher



Curse of Dimensionality

Use Parzen Kernels

Approximate densities for Bayesian classification:

- 2-dimensional kernels

$$p(s_1, s_2) \approx \hat{p}(s_1, s_2) = \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{h} \varphi \left(\frac{x_1(i) - s_1}{h}, \frac{x_2(i) - s_2}{h} \right) \right)$$

- Using independence - approximate 2-dimensional score densities by products of 1-dimensional kernel

$$p(s_1, s_2) = p(s_1)p(s_2) \approx \hat{p}(s_1)\hat{p}(s_2)$$

$$\hat{p}(s_1) = \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{h} \varphi \left(\frac{x_1(i) - s_1}{h} \right) \right)$$

Subscripts *gen* and *imp* have been dropped.

h is estimated using Maximum Likelihood criterion in each case

Estimation Errors

Mean integrated squared error (MISE) is a measure of pdf estimation:

$$MISE(\hat{p}) = E \left(\int_{-\infty}^{\infty} (p - \hat{p})^2(x) dx \right)$$

n = no. of training samples

d = no. of dimensions

k = no. of parameters in derivation;
consider a constant

For d-dimensional kernel estimation of pdf the order of error is:

$$MISE(\hat{p}) \sim n^{-\frac{2k}{2k+d}} \quad \text{W. Hardle "Smoothing techniques with implementation in S"}$$

Known: lower dimension d results in better asymptotic pdf approximation.
In particular, error for approximating 1-dimensional pdf is less than error for approximating 2-dimensional pdf given same no. of training samples

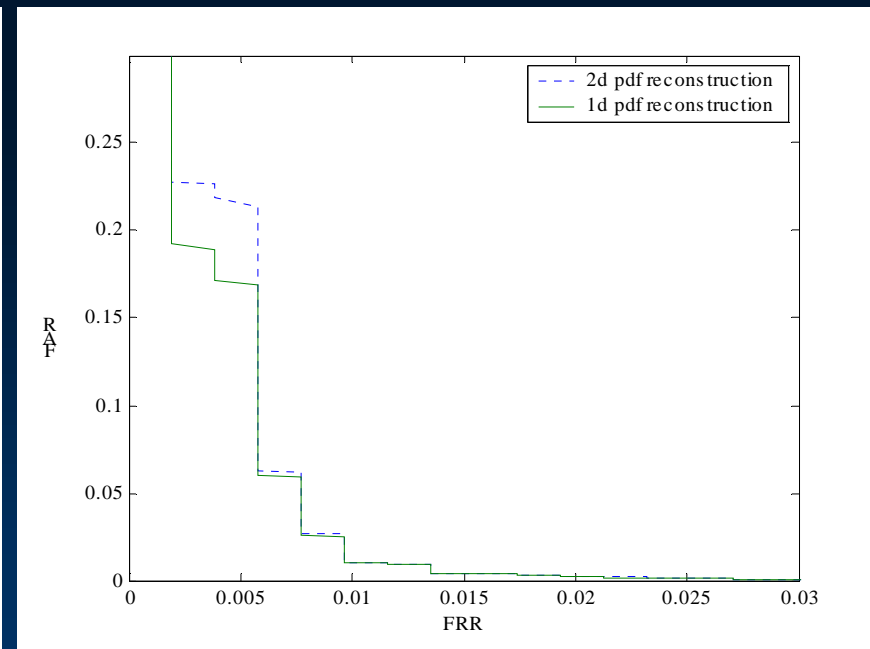
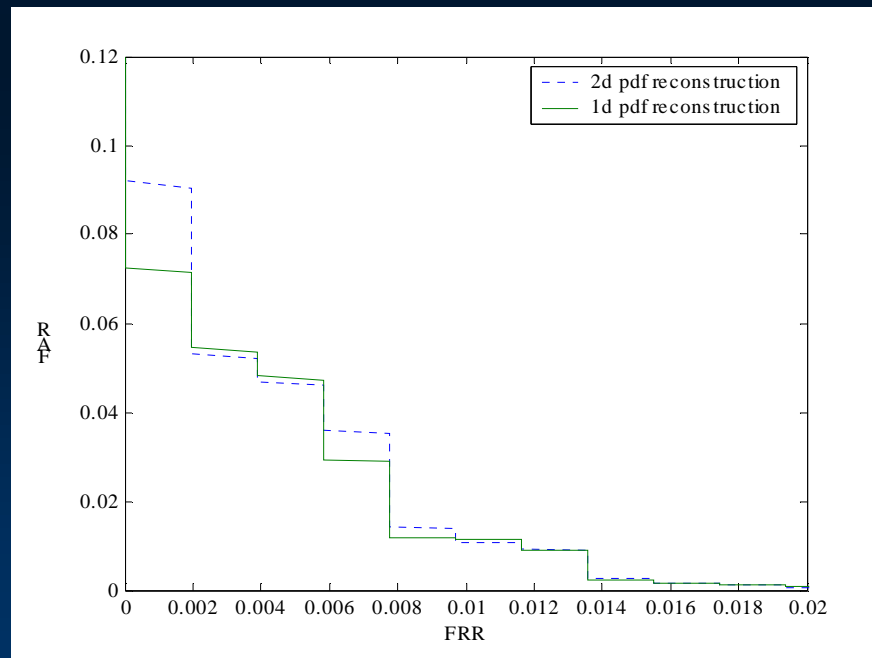
$$n^{-\frac{2k}{2k+1}} < n^{-\frac{2k}{2k+2}}$$

For given n (usually small and fixed) training samples; and constant k

Theorem: Product of approximations has the same order of error as individual 1-dimensional approximations. [Tulyakov 06]

Experiments

- NIST BSSR1 biometric score set
- 517 enrollees and 517 user authentication attempts
- Leave-one-out testing procedure
- Score for test sample (Likelihood Ratio) : $\frac{p(s_1, s_2 | \text{genuine})}{p(s_1, s_2 | \text{impostor})}$
- Can use SVM and other classifiers as well. Not all can produce ROC curve



Summary of Combinations for Verification Tasks

- Can use trainable combination functions instead of fixed rules if classifier is trainable and sufficient training data available
- Can use the score space as feature space
- Likelihood Ratio function is optimal if the probability density functions (PDFs) can be estimated accurately.
- If the individual matchers are independent (as in biometric modalities), the PDFs can be estimated using this information.
- For the same amount of training data (usually limited), estimating lower dimension PDFs is more accurate.