

# Classifiers that improve with use

George Nagy  
DocLab  
Rensselaer Polytechnic Institute

February 19-20, 2004 IEICE-PRMU George Nagy

1

PRMU

## Argument

In-house training sets are **never large enough**, and **never representative enough**. We must therefore augment them with samples from actual (real-time, real-world) OCR operation.

We present some methods to this end.

February 19-20, 2004 IEICE-PRMU George Nagy

2

## Outline

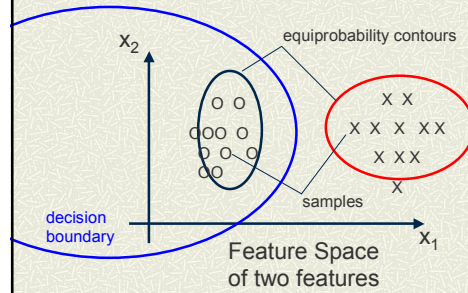
### Non-representative training sets

- Supervised learning (continuing classifier education)
- "Unsupervised" adaptation
  - Self-corrective, Decision-directed, Auto-label
- Symbolic Indirect Correlation (SIC) *new \*\*\**
- Style-constrained classification
- Weakly-constrained data distributions (*new \*\*\**)
- Linguistic context
- Recommendations

February 19-20, 2004 IEICE-PRMU George Nagy

3

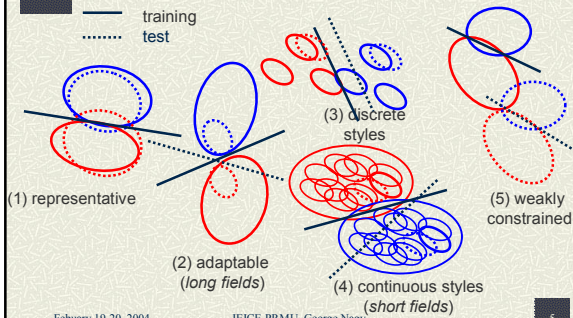
## Representation



February 19-20, 2004 IEICE-PRMU George Nagy

4

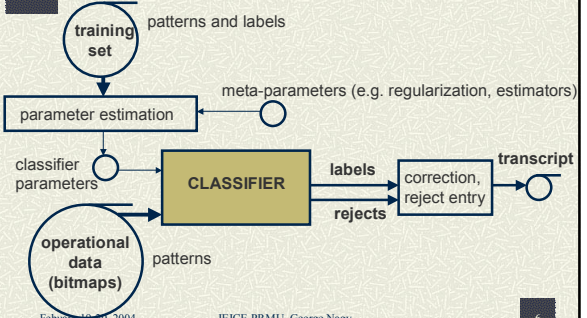
## How representative is the training set?



February 19-20, 2004 IEICE-PRMU George Nagy

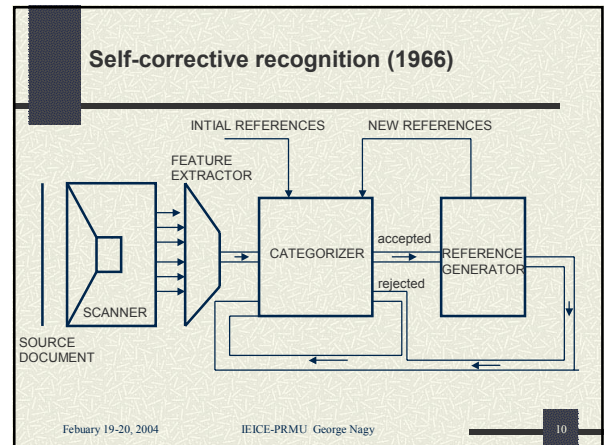
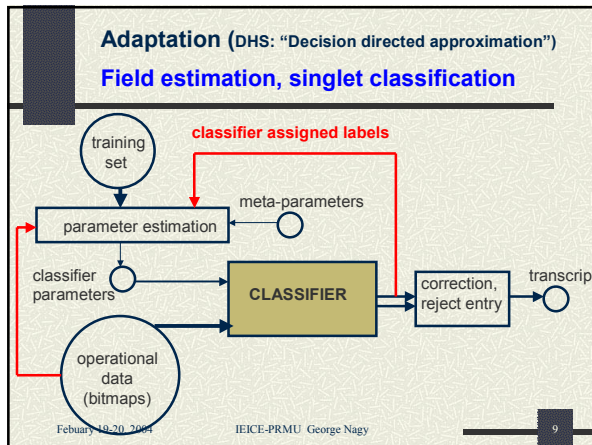
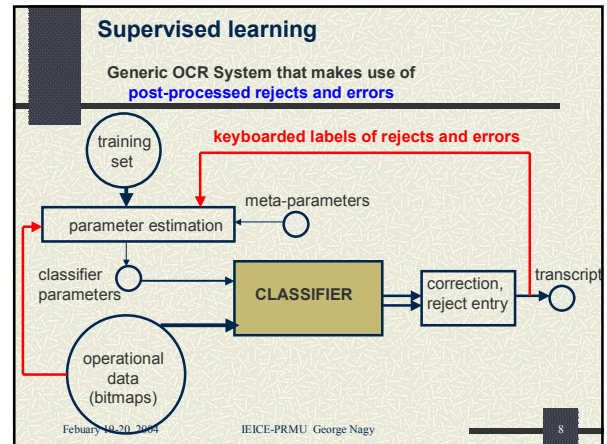
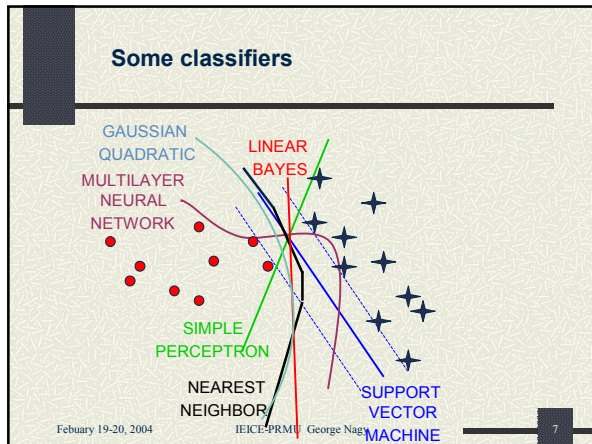
5

## Traditional open-loop OCR System



February 19-20, 2004 IEICE-PRMU George Nagy

6

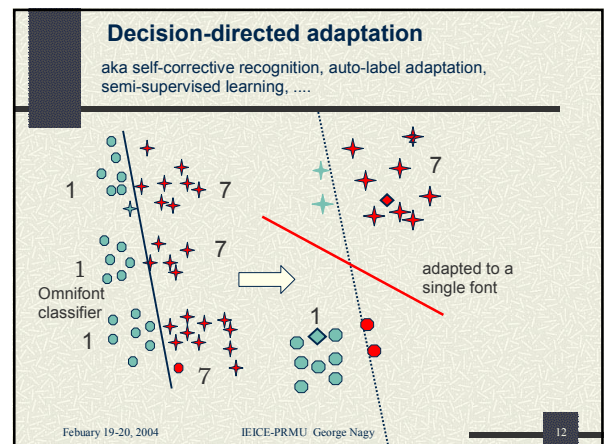


### Results: self-corrective recognition (Shelton & Nagy 1966)

Training set: 9 fonts, 500 characters/font, U/C  
 Test set: 12 fonts, 1500 characters/font, U/C  
 96 n-tuple features, ternary reference vectors

Initial error and reject rates: 3.5% 15.2%  
**After self correction:** 0.7% 3.7%

February 19-20, 2004 IEICE-PRMU George Nagy 11





Query hypothesized as **"contract"**: poor match

February 19-20, 2004 IEICE-PRMU George Nagy 19

Query hypothesized as **"founding"**: good match

February 19-20, 2004 IEICE-PRMU George Nagy 20

Localized Viterbi trellis search

February 19-20, 2004 IEICE-PRMU George Nagy 21

InkLink classification algorithm

1. The **expected location** where the unknown matches the reference words is pre-computed
2. The features matches of the unknown against the reference words are found by **string matching**.
3. The hypothesis that corresponds best to the **expected length and location** of the matches is chosen.

February 19-20, 2004 IEICE-PRMU George Nagy 22

Our most/least favorite writers

Writer-5 An emergency brake cord is located at either end of the car.

Writer-6 An emergency brake cord is located at either end of the car.

Writer-8 An emergency brake cord is located at either end of the car.

Writer-12 How nice of you to totally agree what the founding fathers actually said.

---

Writer-7 An emergency brake cord is located either end of the car.

Writer-9 An emergency brake cord is located at either end of the car.

Writer-10 I like the north rim better than the south rim.

Writer-11 have a laugh, now come hell, but there's a lot to do.

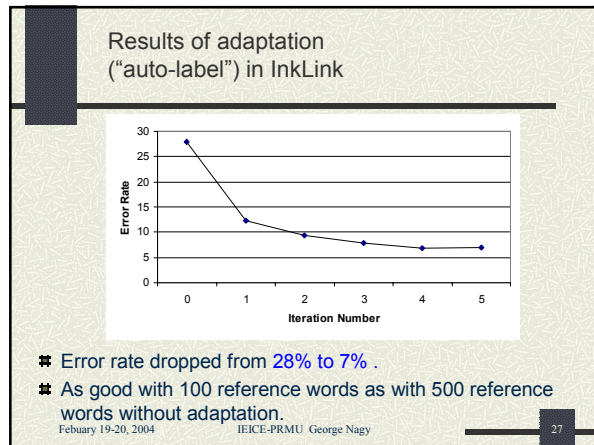
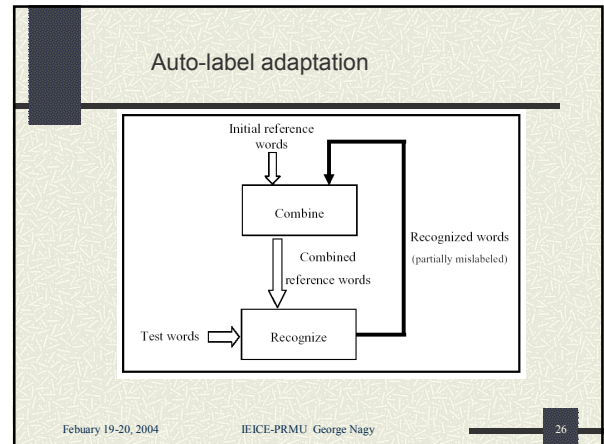
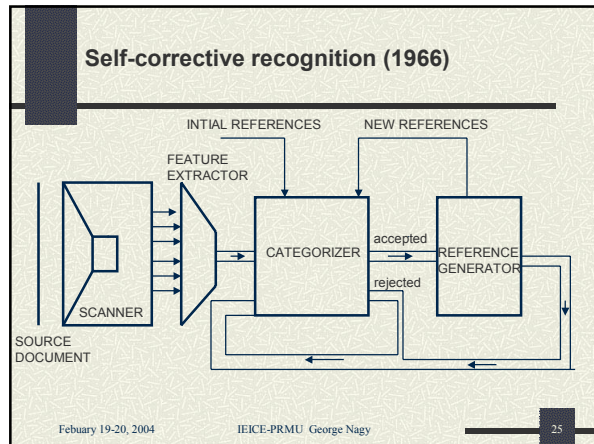
February 19-20, 2004 IEICE-PRMU George Nagy 23

Comparison with external system  
(four writers we like) 100-word lexicons

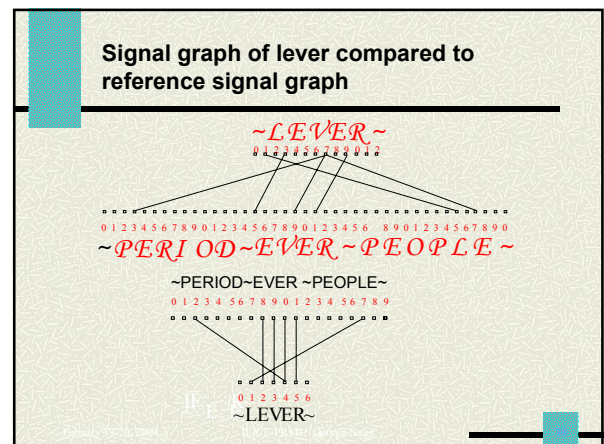
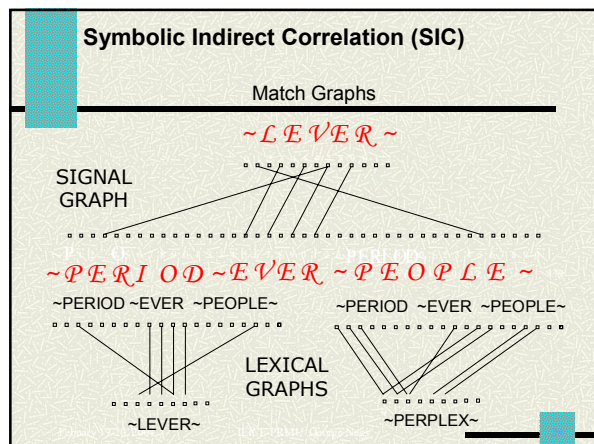
ID \ %	4	5	6	8	Average
%	96.7	95.7	99.7	95.3	96.9
Top-1	96.7	95.7	99.7	95.3	96.9
Top-10	98.3	95.7	99.7	97.7	97.9

ID \ %	4	5	6	8	Average
%	94.0	97.3	94.7	94.0	95.0
Top-1	94.0	97.3	94.7	94.0	95.0
Top-10	99.0	99.0	99.7	99.0	99.2

February 19-20, 2004 IEICE-PRMU George Nagy 24



- ### Outline
- Non-representative training sets
  - Supervised learning (continuing classifier education)
  - "Unsupervised" adaptation
    - Self-corrective, Decision-directed, Auto-label
  - Symbolic Indirect Correlation (SIC)**
  - Style constrained classification
  - Weakly-constrained data distributions
  - Linguistic context
  - Recommendations
- February 19-20, 2004 IEICE-PRMU George Nagy 28



### Matching Match Graphs

Best matching subgraphs preserve Edge Crossings, i.e. find *order-isomorphic subgraphs*

### Signal and lexical graphs (handwriting)

February 19-20, 2004 IEICE-PRMU George Nagy 32

### Results: SIC (simulation only!)

1000 most common words of Brown Corpus

Weighted Noise Model

Probability of Wrong Match:	0.4	0.6	0.8	1.0
% Correct Recognition:	100	100	100	96
Reference String Size:	390	500	600	750

**Adaptation:**  
we can add recognized words to the reference string, as in InkLink

February 19-20, 2004 IEICE-PRMU George Nagy 33

### Outline

- Non-representative training sets
- Supervised learning (continuing classifier education)
- "Unsupervised" adaptation
  - Self-corrective, Decision-directed, Auto-label
- Symbolic Indirect Correlation
- Style constrained classification**
- Weakly-constrained data distributions
- Linguistic context
- Recommendations

February 19-20, 2004 IEICE-PRMU George Nagy 34

### Style consistency: Field estimation, field classification

February 19-20, 2004 IEICE-PRMU George Nagy 35

### Single-class and multi-class style

SINGLE CLASS STYLE	MULTI-CLASS STYLE
Source 1: 29/05/1925	25/07/1922
Source 2: 15/05/1990	05/05/1925
Source 3: 21/06/1943	02/06/1943
Source 4: 05 /29/1945	02/25/1942

Styles are induced in a collection of documents by multiple sources\*.

\* fonts, printers, scanners, writers, speakers, microphones, ...

February 19-20, 2004 IEICE-PRMU George Nagy 36

### Multi-mode parametric classifier for single-class style constraints

February 19-20, 2004 IEICE-PRMU George Nagy 37

### Field classifier for strong style constraints

February 19-20, 2004 IEICE-PRMU George Nagy 38

### Field-trained classification versus style-constrained classification

Field Length = 4

Training set for field classification	Training set for style classification
0000	00
0001	01
0010	02
	98
	99
	(10 <sup>2</sup> classes)
9998	
9999	
(10 <sup>4</sup> classes)	

classifier parameters for longer field length computed from pair parameters (because Gaussian variables defined completely by covariance)

February 19-20, 2004 IEICE-PRMU George Nagy 39

### Results: style-constrained classification - short fields (Harsha V.)

Continuous style-constrained classifier, trained on ~ 17,000 characters and tested on ~17,000 characters.  
25 top principal component "Hitachi" blurred directional features.

Test data	Field error rate (%)			
	Field length: L=2		L=5	
	w/o style	with style	w/o style	with style
SD3	1.4	1.3	3.0	2.5
SD7	2.7	2.4	5.3	4.5

February 19-20, 2004 IEICE-PRMU George Nagy 40

### Outline

- Non-representative training sets
- Supervised learning (continuing classifier education)
- "Unsupervised" adaptation
  - Self-corrective, Decision-directed, Auto-label
- Symbolic Indirect Correlation
- Style constrained classification
- Weakly-constrained data distributions**
- Linguistic context
- Recommendations

February 19-20, 2004 IEICE-PRMU George Nagy 41

### Weakly-constrained data

given  $p(x)$ , find  $p(y)$ , where  $y=g(x)$

3 classes, 4 multi-class styles

February 19-20, 2004 IEICE-PRMU George Nagy 42





## Style context versus Language context

Two digits in an isogenous field: ..... 5 6 .....

with feature vectors:  $\mathbf{x}_i \mathbf{x}_j$

and class labels:  $C_i C_j$

Style context:  $P(\mathbf{x}_i, \mathbf{x}_j | C_i=5, C_j=6) \neq P(\mathbf{x}_i | C_i=5) P(\mathbf{x}_j | C_j=6)$

Language Context:  $P(C_i=5, C_j=6) \neq P(C_i=5) P(C_j=6)$

February 19-20, 2004

IEICE-PRMU George Nagy

49

## Recommendations for OCR systems that improve with use

Never let the machine rest:  
design it so that it puts every coffee-break to good use.

Don't throw away edits (corrected labels): use them.

Classify style-consistent fields, not characters:  
*adapt* on long fields, exploit multi-class *style* in short fields.

Use order rather than position.

Let the machine guess: lazy decisions.

Make use of all possible context:  
language, shape, layout, and function.

Every document can be read by the right reader.

February 19-20, 2004

IEICE-PRMU George Nagy

50

## The End

I am grateful to Hitachi CRL for technical and financial support since 1992. Special thanks to Drs. Fujisawa, Liu, Sako, and Shima, and to Messrs. Koga, Marukawa, and Mine.

I also learned much from my former and present students Jung, Sarkar, Veeramachaneni, & El-Nasan.

But any misconceptions in this presentation are entirely my own!

Thank you!

February 19-20, 2004

IEICE-PRMU George Nagy

51

## Non-supervised classification is always semi-supervised classification

*Only the supervisory constraints are hidden!*

Clustering: cardinality, diameter, population, or sequence of clusters.

ML / EM : form of distribution, sigma, number of mixture components.

Trees, neural networks: validation set.

February 19-20, 2004

IEICE-PRMU George Nagy

52

## InkLink classification algorithm

$$P(c_i | l_{ij}) = \frac{P(l_{ij} | c_i) P(c_i)}{P(l_{ij})}$$

$$c^* = \arg \max_c \prod_i \frac{P(l_{ij} | \hat{l}_{ij}, M_{ij} = 1)}{P(l_{ij} | \hat{l}_{ij}, M_{ij} = 0)}$$

$l_{ij}$  : Observed feature match length

$\hat{l}_{ij}$  : Predicted feature match length

$M_{ij}$  : Match

Four classifiers, based on different features sets, are combined by Borda Count

February 19-20, 2004

IEICE-PRMU George Nagy

53

## A letter from George Washington



Feature-string matching:

"complete" vs. "are to be made to the aid de camp."

February 19-20, 2004

IEICE-PRMU George Nagy

54

## CONFERENCE REPORT ON ICDAR '07

George Nagy  
Professor Emeritus (??), Rensselaer Polytechnic Institute

### OUTLINE:

VENUE  
INVITED SPEAKERS  
AWARDS  
SESSION TOPICS  
NOTABLE CONTRIBUTIONS  
LARGE-SCALE PROJECTS  
OUTLOOK

(Slides from my closing talk at ICDAR 1997 in Ulm, Germany)

February 19-20, 2004

IEICE-PRMU George Nagy

55

## INVITED SPEAKERS

(all second generation)

**Prof. Tin Kam Ho**  
Web-wide Voting Network for Weak Classifiers  
Exploratory data analysis at Lucent Bell Labs

**Dr. Tapas Kanungo**  
Multimedia Document Devastation Models  
DIA at IBM Almaden Research Center

**Dr. Rainer Hoch**  
High-speed Context Switches  
Just moved from SAP to the Mannheim Berufsakademie

**Dr. Abdelwahab Zramdini**  
Web Typography  
Top secret research at a bank in Geneva

**Prof. Omid Kia**  
Processing Compressed Multimedia Documents  
Founded a company, then disappeared!

February 19-20, 2004

IEICE-PRMU George Nagy

56

## III. Devices

### Storage

Docupin, DocuPinCushion (DPC)  
750 MB/DP, 18GB/DPC  
Not addressable

... イベントのお知らせ "Memory Stick Xmas Magic". 楽しくお得な各種  
キャンペーンをご紹介します!! ... 最新情報誌 Memory Stick Update. ...  
SANDISK UNVEILS THE HIGHEST CAPACITY MEMORY STICK PRO STORAGE  
CARD IN THE WORLD—2 GIGABYTES.



February 19-20, 2004

IEICE-PRMU George Nagy

57

## E-Scap

Cholesteric polymer-dispersed hybrid substrate  
10 micron thick (100 sheets per cm)  
Ten megacell array provides 300 dpi resolution  
10\*\*6 read-write cycles  
Ecologically benign (silicon, not carbon, based)  
No color so far



February 19-20, 2004

IEICE-PRMU George Nagy

58

## Radiotablet

Lambertian high-contrast reflective display  
Satellite cyberband channel  
Pentameric (soft, 4mm thick elastomere) vary-keyboard  
100 MB flash-cache  
Docupin port (USB)



### Tablet PCs

Toshiba's Convertible Tablet PC paves the way ...



February 19-20, 2004

IEICE-PRMU George Nagy

59

## VI. Oral Documents

Speech-to-print  
Print-to-speech  
Convergence of speech and document image processing -  
began in 1990's with adoption of HMM in OCR

### AUDIO DOCUMENTS

So far, primarily for entertainment (music)  
(almost everything is primarily for entertainment!)  
Some natural sounds (bird calls, frogs, dolphins)  
Oral history (eye-witness accounts)  
Personal note-taking (key-chain recorders)  
Instructional audios ("prompter")  
???

February 19-20, 2004

IEICE-PRMU George Nagy

60

## handwriting

Now used mainly for *self-communication*

letters already replaced by email

forms being replaced by e-forms

Children learn to type in pre-school!

Therefore personal electronic ink recognition  
is becoming more important than  
scanned handwriting recognition.

February 19-20, 2004

IEICE-PRMU George Nagy

61

## Other questions?

February 19-20, 2004

IEICE-PRMU George Nagy

62