

Semi-supervised Learning for Handwriting Recognition

Gregory R. Ball and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition
University at Buffalo, USA

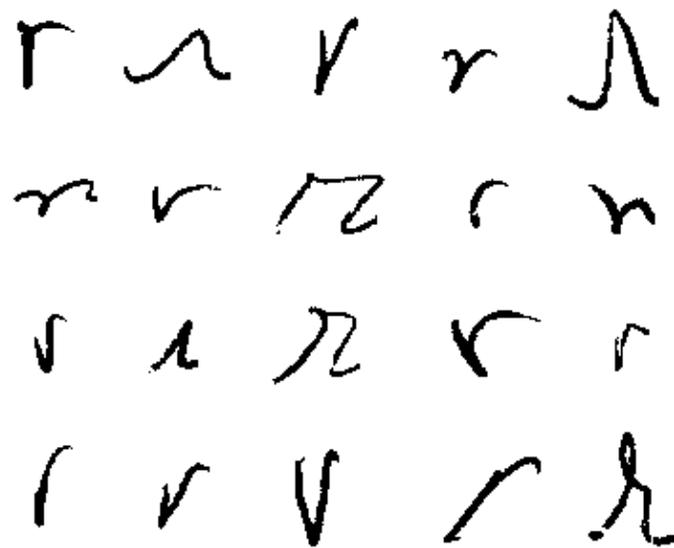
ICDAR 2009, Barcelona, Spain

Overview

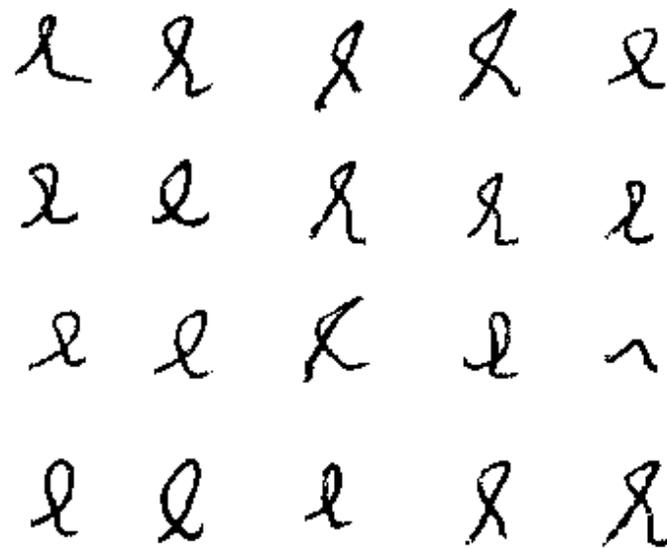
- Writer Adaptation / Semi-supervised Learning
 - Determining uniqueness of writing
- Algorithm
 - Feature set and recognition
 - Datasets (English and Arabic)
 - Results
- Concluding Remarks

Writer Adaptation

- Every writing style not captured in *training* set
- Task of *learning* particular writer's style in recognition



Letter "r" many authors



Letter "r" one writer

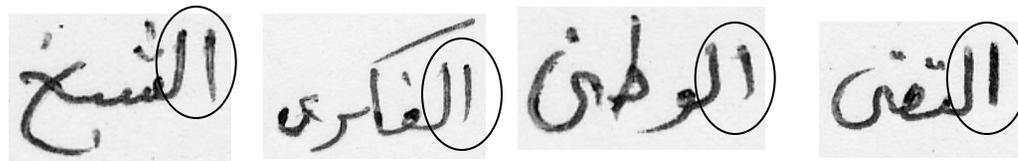
Writer Adaptation

- Style-dependencies occur in every language

- Arabic example

- Writer A

- Training



- Writer B

- Testing



- Writer A has standard *alef-lam* combinations (most prototypes like this)
- Learning writer B style can increase recognition rate

Writer Adaptation Literature

- On-line handwriting recognition
 - Personalized systems by using supervised learning (enrollment adaptation)
 - Haluptzok, Revow, Abdulkader (*IWFHR 06*)
 - Subrahmonia, Nathan, Perrone (*ICASSP 96*)
 - Senior and Nathan (*ICASSP 97*)
- Off-line recognition
 - Personalized systems
 - Bozinovic and Srihari (*IEEE PAMI, 93*)
 - Creating prototype set exclusively from author samples
 - Vinciarelli and Bengio (*PRL 02*)
 - Recognition output as training data
 - Nosary et al. (*IWFHR 02*)

Semi-Supervised Learning

- Motivation
 - Mixture of labeled/unlabeled data available
- Use of Training Data
 - Labeled data are known exemplars
 - Unlabeled data used to refine learning
- Value of Approach
 - Exploit large quantities of training data while avoiding expense in labeling

Semi-supervised Learning Literature

- Semi-supervised learning
 - S. Fralick. Learning to recognize patterns without a teacher. IEEE Trans Inform Theory, Jan 1967.
 - O. Chapelle, B. Scholkopf, and A. Zien, eds. Semi-Supervised Learning. MIT Press, 2006.
 - X. Zhu. Semi-Supervised Learning with Graphs. PhD thesis, CMU, 2005.
- Word sense disambiguation
 - D. Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. Proc ACL 1995.
- Distinguish subjective/objective sentences
 - E. Riloff, J. Wiebe, and T. Wilson. Learning subjective nouns using extraction pattern bootstrapping. Proc 7th Conf Natural Language Learning, 2003.
- Emotions in spoken dialogue
 - B. Maeireizo, D. Litman, and R. Hwa. Co-training for predicting emotions with spoken dialogue data. In Proc ACL 2004.

Large-scale labeled training in handwriting

- Unavailability (for many scripts except English)
- Diversity of writing styles leads to over-generalized models
 - for a single writer, generalized model is not best
 - Can be overcome with manual training, but such supervised learning can be expensive

Semi-Supervised Learning in Handwriting

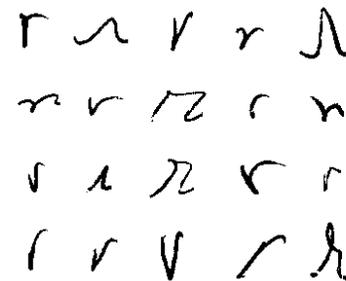
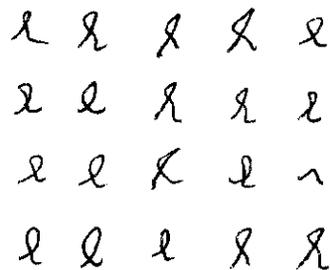
- More cost effective solution
- Generalized writing model generated by in advance
 - Utilizing large dataset of manually truthed documents
- Runtime sample treated as unlabeled set
- Learning takes place concurrently with recognition task
 - Need some changes to standard approaches, such as utilizing self-learning

Semi-Supervised Approach

- Use writer independent system to select samples within input document for use as unlabeled training data
- Combining HMM recognition engine with writer adaptation
- Explored how to integrate author's sampled writing into a writer dependent model.
- Data sets to evaluate performance
 - English, tested on the CEDAR letter database
 - Arabic, tested on a CEDAR Arabic dataset

Detecting Writing Uniqueness

- Utilize unusual characters that are still internally consistent.
- Mimic human reading with non-standard characters
- Generalized models fail to capture this idea since such characters not weighted heavily when creating the model.
 - Writer Verification ignored them as “bad prototypes”
 - Supervised models need expensive “enrollment” step



Letter “r” written by a single writer CEDAR

Letter “r” written by many writers

Character Feature Set



2 Global Features: Aspect ratio
Stroke ratio

72 Local Features

$$F(i,j) = \frac{s(i,j)}{N(i)*S(j)} \quad \begin{matrix} i = 1,9 \\ j = 0,7 \end{matrix}$$

$s(i,j)$ - no. of components with slope j in sub image i

$N(i)$ - no. of components in sub image i

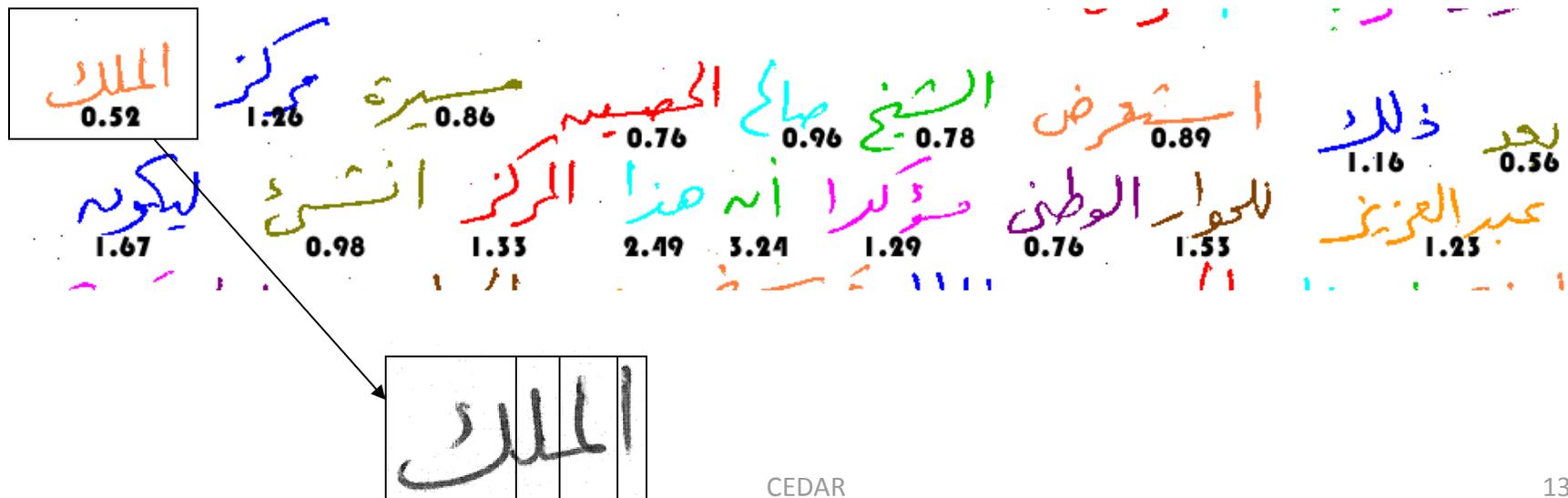
$$S(j) = \max \frac{s(i,j)}{N(i)}$$

0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.076923 0.000000 0.529915 0.993590 0.387363 0.000000 0.596154 1.000000	0.175000 0.962500 0.688889 0.129167 0.419643 1.000000 0.387500 0.000000
1.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.727273 1.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.337662	0.000000 0.000000 1.000000 1.000000 0.108295 0.000000 1.000000 0.599078
0.207547 0.207547 0.194969 0.779874 0.760108 0.215633 0.350943 0.630728	0.127660 0.351064 0.073286 0.879433 1.000000 0.607903 0.263830 0.474164	0.111111 0.407407 0.637860 0.765432 0.870370 0.000000 0.114815 0.687831

WMR Feature Vector (74 values)

Adaptation Process

- Entire page recognized
- Words with highest recognition confidence isolated
- Features extracted from characters
- Model adjusted (new features added or combined)
- Process repeated



Algorithm

- Use writer independent model to recognize all words W in document D . Associate each word $w \in W$ with a score w_s .
- Rank all words by their corresponding w_s scores.
- Select top- n words
- Incorporate features into writer-independent model
 - Split words into sub-characters
 - Combine into characters according to truth value
 - Split characters
 - Extract features
 - Combine features into feature set (add, combine, replace, etc. feature vectors in independent model)
- Iterate (repeat starting at first step) until no significant increase in overall confidence of D (e.g. $\sum w_s$).

CEDAR Letter Dataset

- 156 words including all characters (letters and numerals), punctuations and distinctive letter and numeral combinations (ff, tt, oo, 00)
- About 1, 500 individuals copied CEDAR letter three times
- Samples scanned:300 dpi, 8-bit grayscale.

From Nov 10, 1979
Jim Elder
829 Loop Street, Apt 300
Allentown, New York 14707

To
Dr. Bob Grant
602 Greenberry Parkway
Omae, West Virginia 25638

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Zack.

It all started about six months ago while attending the "Rubey" Jazz Concert. Organizing such an event is no picnic, and as President of the Alumni Association, a co-sponsor of the event, Kate was overworked. But she enjoyed her job, and did what was required of her with great zeal and enthusiasm.

However, the extra hours affected her health; halfway through the show she passed out. We rushed her to the hospital, and several questions, x-rays and blood tests later, were told it was just exhaustion.

Kate's been in very bad health since. Could you kindly take a look at the results and give us your opinion?

Thank you!
Jim

Writer Independent Recognition

- Writer independent model
 - HMM using a lexicon (*IWFHR 06*)
 - Trained on IFN/ENIT dataset for Arabic, internal datasets for English
- Several trained sets attempted (5-100 clusters), five cluster version used for most recognition tasks
- Approximately 200 character classes (different classes for different positions, compound classes)

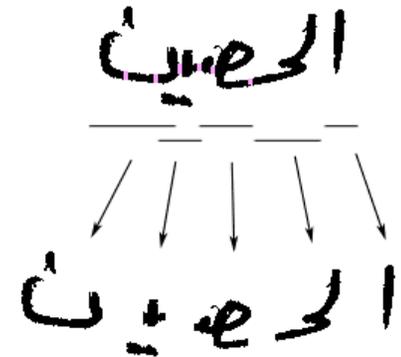
CEDAR Arabic Page Dataset

- Used for testing
- Ten one page articles written by ten authors
- Manually truthed to word level



Arabic Recognition Engine

- Character-based algorithm
- Based on word model recognition feature set
- Trained on IFN/INET dataset
- Several trained sets attempted (5-100 clusters), five cluster version used for most recognition tasks
- Approximately 200 character classes (different classes for different positions, compound classes)

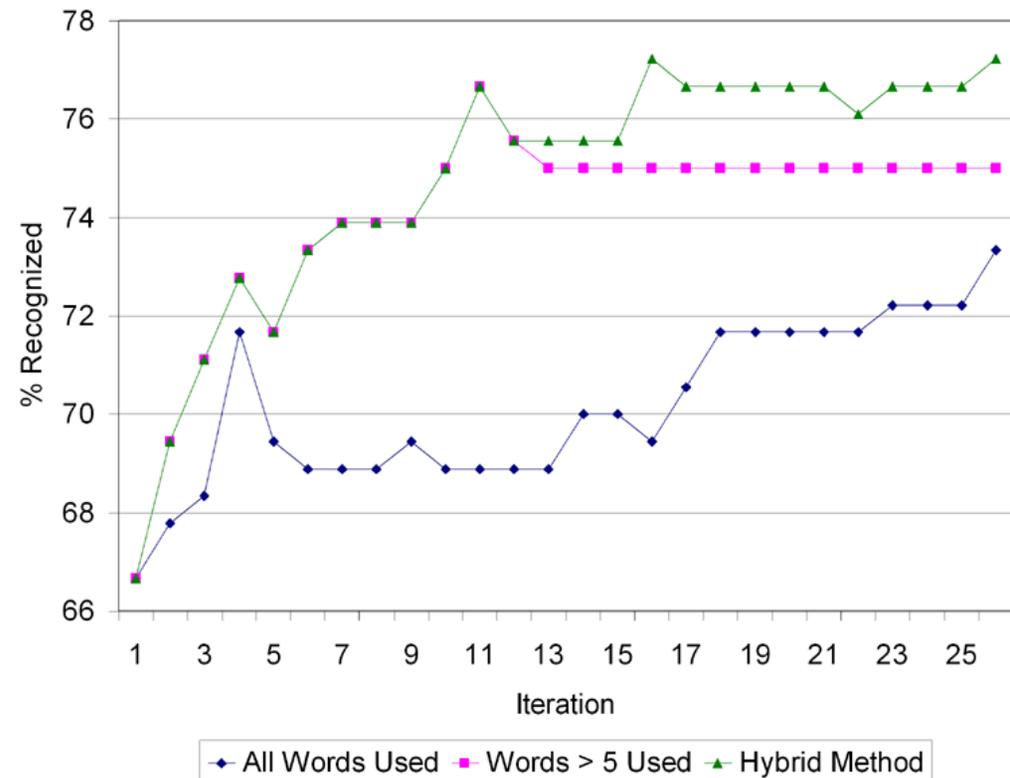


Improved Confidence Metric

- *Short* incorrectly recognized words are likely to still have high confidence scores.
- Less likely with *long* words since very different words are unlikely to have ambiguous segmentation points
- Poorly recognized short words can have a detrimental effect on the adaptation process
- Strategy is to use longer words first

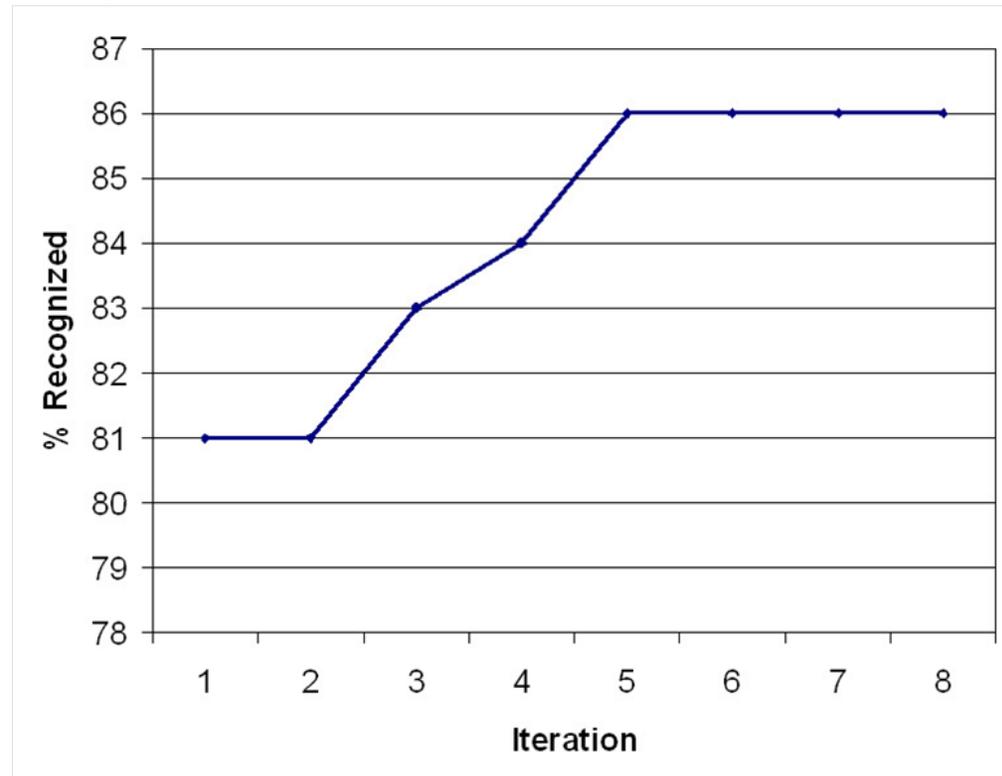
Arabic Adaptation Results

- Adapting to longer words first provides largest improvement
- Each iteration adds five words
- Large initial increase that levels off
- Results in about 10% recognition improvement



English Adaptation Results

- Significant improvement
- Less magnitude due to smaller character set cardinality and higher starting recognition rate



Results Summary

- For Arabic, recognition moved from 67-77%
 - After basic adaptation overall recognition rate was 73%.
 - Focusing on longer words improved from 70% recognition to about 92% recognition rate after adaptation.
 - Hybrid method is achieved a rate of about 77%
- English had similar results, overall moving from 81-86%

Adaptation Costs

- Iterative model (worst case) reproduces much recognition effort
- E.g., good increase after ten iterations—naïve algorithm bottleneck of ten full recognition attempts at D .
- Some improvement by only rescoreing words affected by new feature set
 - (Often can) exclude words used for adaptation

Concluding Remarks/Future Directions

- Significant increase in performance based on a small sample size (several words or more)
 - Writer independent model adaptation
 - Future direction: “ground up” model
 - Additional sample size requirements
 - Potentially higher performance increase
- Applicable to many other techniques
 - Perform adaptation as preprocessing step
 - Convert algorithms to iterative versions

Thank You!

- For further information:

srihari@cedar.buffalo.edu

grball@cedar.buffalo.edu