Artificial Intelligence at the Post Office and the Police Department

Sargur Srihari

University at Buffalo The State University of New York

From Mail Sorting to CSI





Three Part Story

- 1. Postal Address Interpretation
- 2. CSI
- 3. Science of AI-PR-ML

Handwritten Address Interpretation



Postnet Bar Code representing Delivery Point

*Among first handwritten addresses sorted automatically (Oct 1996)

Handwriting Recognition

Characters





68 86 54 22 80 20 60 30 80 89 26 00 20 53

80 50 21 08 80 80 00

33 (\$ 59 44 00 43

20 80 84 99 00 80 77 00 22 80 20 Letter Pairs (state abbr)





H Held In



MN Muhio My Pa 1

Ut Corr.

Words

Mining Postal Directories



Context Strategy

Street address



USPS Handwritten Address Encode Rate

1997: 30% finalization

2001: 75% finalization http://www.usps.com/history/cs01/c2f.htm 2004: 88% finalization

From 1996 to 2004, the encode rate for RCRs increased from 35 percent to 90 percent, reducing the need for manual matching at RECs from 24 billion pieces per year to 6 billion annually. As a result, USPS has been able to reduce the number of RECs in the national network from a high of 55 in 1998 to only 15 today.

AI Technologies

- Character recognition ← image analysis/pattern recognition
- Cursive word recognition \leftarrow HMM's
- Street, City identification ← constraint satisfaction
- Address interpretation ← use of background knowledge of legal postal addresses
- Case study in:
 - Knowledge-driven image interpretation
 - Machine learning

Knowledge-Guided Interpretation

Knowledge base of

Street address



Lessons Learned

- 1. Context overcomes handwriting ambiguity
- 2. Statistical methods allowed discovery of underlying relationships in designing recognizers
- 3. Machine performance exceeds human performance in a skill not highly specialized

Deployment and Application

USPS handles over 70 million handwritten addresses daily

- 83% of these are now automatically sorted
- Machine automatically adds 9-digit zip code to each piece
- 13 pieces processed / second

Initial development:

- SUNY Buffalo, 15 years of research & development

Ongoing refinement:

- Lockheed Martin, Siemens

Impact and Performance Metrics

Coverage:

 – Over 55 million handwritten addresses sorted daily (83% of all handwritten)

Accuracy: – above 98%

Savings:

- \$100M saved during first year deployed (1997)
- Larger annual savings today, with increased accuracy

USPS Follow-up

- Continuing improvements in accuracy
- Adoption by foreign countries

 Australia, UK





• Writer identification



Landmarks in Postal Research at CEDAR

- 1980: NSF grant on "Contextual OCR Algorithms" to Srihari, P.I.
- 1983: First USPS contract awarded
- 1993: Group designated by USPS as CEDAR
- 1997: First deployment of HWAI by USPS in Tampa
- 1999: Australia Post Deployment
- 2000: UK Royal Mail Deployment
- 2008: Flats, Evaluations

Line Segmentation

Touching lines

Nor 10, 1999
From
Jim Eldar
829 Loop Scheet, Apt 300
Allentonn, Non Jork 14707
To D 7 1 0 1
Ur. Dob Grant
Oner West Troining 25638
We prove poloured to an I a Xeno Choin at the University
Medical Center. This is regarding my friend, Kate Zack.
It all started aroud six months ago while attending the
" Rubeg" Jozy Concort. Organizing such an even and
a - sponsor of the event, tate was overworksed But she
enjoyed her joh, and did what was beguined of
Great Zeal and entrustant.
Ih and the autor Reported for to althe halfway thinghe The
Show she passed out. We rushed her to the holystal)
Several greatures 1x - rays in practice of the
told it was go in findly
Kate's been in very had health smae, could of prinions
take a look at the verie - 0.
Thouk you!
· · · · ·

Overlapping lines

a N

Line and Word Segmentation

Allenbourn, New Jork 14707 Dr. Bob grant 602 queensberry Parkway Omar, west urging 25638 Ave were referred to you by xena When at the uni-versity Reducal Center. This is regarding my firend, hate Zack It all plotted bit months ago while the user the user a planing the user is no punie as planet is no punie as planet is no punie "Rubeo" H of the did what was required of her with great peat However, the extra hours officed her health; halfway Anough the blue she passed traf. We rushed by to the hospitol, and bevere suchions, know blood toot lober, were told it was full entruition. Kate's beenin very bad health Sine. Could 11. 0 0 10 0.01.04 0.35

Handwriting Thesis Defense, Summer 2008, Lund, Sweden



Legal Motivation

- Supreme Court Ruling regarding expert testimony
 - Daubert vs. Merrell Dow (1993) technique needs to have a *scientific* basis:
 (i) testing, (ii) error rates, (iii) peer review and (iv) acceptability
- Scientific Basis
 - Philosophy of Science: Scientific Theory Falsifiability- Popper
- Legal Clarifications
 - GE *vs*. Joiner (1997)

reliable scientific evidence must be admitted. Unreliable excluded. "Weight of evidence" methodology acceptable

• Kumho vs.Carmichael (1999)

technique produces consistent result whether scientific, technical or other specialized knowledge

Importance of Handwriting in Forensics

- Bureau of Justice Statistics (2002)
 - Among 50 largest publicly funded crime labs
 - 57% perform QD function
 - 5,231 cases requested
 - 1,079 backlogged at year end
- Significantly larger case load internationally
- Handwriting is common in QD case work



Computer Of The summer Their be back Jame 2. 3 may 4 for the (Radial heis Materials) . A down to see you there Many Marks again Kee see Jon is the face of leas





Questioned Document Examination (QDE) is an Expertise - Based on Human Cognitive Skills

Same Writer?

Nov. 10, 1999 Movember 10, 1994 From Jim Elder 829 Loop Street, Apt. 300 Jim Elder Allentaun, New York 14707 829 Loop Street, Apt. 200 То Allenton New Yok 14707 pr. Bob Grant 602 Queensberry Parkway To Ar. Brs Grant Omar, West Virginia 25638 602 Queensberry Parking We were referred to you by Xena Cohen at the University Omar , west Kingina 25638 Medical Center. This is regarding my friend, Kate Zack. It all started around six months ago while attending the We were reperred to you by Xena Cohen at the "Rubeq" Jazz concert. Organizing such an event is no picnic, and as president of the Alumni Association, a co-sponsor University Medical Center. This is reparding my friend Kate Zack of the event, take was averworked. But she enjoyed her job, It all started around six months ago while altendie and did what was required of her with great zeal and enthusiasm, the "Ruber" Jz; consent. Organizing such an event However, the extra hours affected her health; halfway through is no picnic, and as Prevident of the Alumi Association the show she passed out. We rushed her to the hospital, a conspirer of the event, Kale was overworked. But and several questions, x-rays and blood tests later, were told she apped her job, and did what was required of her it was just exhaustion. Kaks been in very bad health since. Could you kindly take with great zeel and enthis arm. a look at the results and give us your opinion? However, the extra have affected her health helfing Thank You! Things the she she posed out. We reached her to The prospiral, and several question, X-rays and Jim blood texts later, were tild it in just exhaustin

Information Processing Model of QDE

1. Computational Theory

- Goal of Computation: Verification/Identification
- Logic of Strategy: Feature Space/Similarity Space

2. Representation and Algorithms

Handwriting Elements
 Discriminating Power

3. Implementation

- Software System
- Define a standard: Strength of Evidence
- Evaluate

Repeat

CEDARFOX: A Software System for the Comparison of Extended Handwriting and Signatures



Handwritten Document Analysis System

Writer Identification

CEDAR-FOX - 0001au.png Questioned Document	
Questioned Document	
From Jan Elder 237 Lap Shaed, ppt 200 Alkantenin, New York 14707 To To To To To To To To To To	Ranked Document List

Probability of Random Correspondence (PRC)

- Birthday Paradox
 - Discrete variable
- Human Height
 - Continuous variable, needs tolerance to be specified
- Fingerprints
 - Three variables, needs tolerance

PRC in Birthday Problem

- PRC
 - probability that two people have *P* the same birthday

$$PRC = \sum_{1}^{365} \left(\frac{1}{365}\right)^2 = \frac{1}{365} = 0.0027$$

- General PRC
 - probability that in a group of people (n), some pair of them have the same birthday

$p(n) = 1 - (1 - PRC)^{\binom{n}{2}}$	
with $n = 2$, $p(2) = PRC$	

n	p(n)		
2	0.0027		
5	0.0271		
10	0.1169		
20	0.4114		
40	0.8912		
80	0.9999		
120	0.99999997		
370	1		
510	1		

PRC in Birthday Problem

- PRC
 - two people have same birthday
- General PRC (birthday paradox)
 - some pair among *n* have same birthday

$$PRC = \sum_{1}^{365} \left(\frac{1}{365}\right)^2 = \frac{1}{365} = 0.0027$$

$$p(n) = 1 - (1 - PRC)^{\binom{n}{2}}$$

with $n = 2, p(2) = PRC$

- Specific PRC
 - probability that given a person's birthday (b) in a group of (n) people, at least one person shares the same birthday (b) among other (n -1) persons

$$p(n,b) = 1 - (1 - p(b))^{n-1} = 1 - (1 - \frac{1}{365})^{n-1}$$

Human Height: PRC



	Female	Male
Mean	5'3"	5'8"
Standard Deviation	11.1"	3.3"

$$p_{\varepsilon} = \int_{-\infty}^{\infty} \left(\int_{\alpha-\varepsilon}^{\alpha+\varepsilon} P(h \mid \mu, \delta) dh \right)^2 da$$

$$P(h \mid \mu, \delta) \text{ is the probabilistic generative model}$$

$$\varepsilon \qquad \text{ is the tolerance}$$

PRC for female height with $\varepsilon = 0.1''$ is 0.0025 PRC for male height is 0.0085

Human Height: General and Specific $p(h \mid \mu, \delta) = \frac{1}{\sqrt{2\pi\delta}} e^{-\frac{(h-\mu)^2}{2\delta^2}}$

General PRC

Specific PRC

$$p(n) = 1 - (1 - PRC)^{(\xi)}$$

(n)

 $p(h,n) = 1 - (1 - p(h \mid \mu, \delta))^{n-1}$

	p(n)		
n	Female	Male	
2	0.0025	0.085	
5	0.0247	0.0818	
10	0.1065	0.3190	
20	0.3785	0.8025	
40	0.8581	0.9987	
80	0.9996	0.9999999992	
120	0.99999	1	
370	0.9999999999	1	
	99		
510	1	1	

 $P(h \mid \mu, \delta)$ is probability person has height h

	p(n,h)		
n	Female (57 inches)	Male (68 inches)	
2	0.0030	0.0112	
5	0.0248	0.0825	
40	0.1098	0.3551	
80	0.2100	0.5888	
140	0.3395	0.7906	
200	0.4477	0.8934	
400	0.6959	0.9888	
800	0.9078	0.9989	
1000	0.9492	0.9998	
2000	0.9974	0.9999994	

Fingerprints: General PRC

• General PRC in 100,000 fingerprints

THE GENERAL PRCs with varying M and W in n = 100000 fingerprints

		Minutiae and Ridge Information	only Minutiae
t	m	p(n)	p(n)
46	46	6.7329×10^{-49}	5.0680×10^{-34}
	36	5.9220×10^{-27}	1.3336×10^{-15}
	26	1.7579×10^{-11}	1.1842×10^{-3}
	16	1	1
	6	1	1
36	36	7.5000×10^{-35}	1.5551×10^{-22}
	26	1.5202×10^{-14}	5.3414×10^{-6}
	16	3.4852×10^{-1}	1
	6	1	1
26	26	5.9297×10^{-20}	2.3517×10^{-11}
	16	1.8829×10^{-2}	1
	6	1	1
16	16	4.460×10^{-8}	1.3464×10^{-1}
	6	1	1

$$p(n) = 1 - (1 - PRC)^{\binom{n}{2}}$$

(t) the number of minutiae in template fingerprint.

(m) the number of matched minutiae(n) the number of fingerprints.

Shoeprint Evidence in Crime Scenes

- Most commonly found evidence at crime scenes
- Can be used for identification of suspects by narrowing down search space
 - By elimination of shoe type, by matching it against a set of known shoeprints (impressions of many different types of shoes on a chemical surface)
- Manual identification is laborious
 - Hence seek automated methods

Types of Latent Shoe Prints Positive print Negative print





Databases

- Latent Print database
 - Pairs of images
 - latent and chemical footwear prints
- Shoe Patterns Database

Terminologies

- Latent Print: The shoe print which are left by the criminals at the crime scene (also known as unkown print)
- Chemical Print: The shoe print taken from a known sole on a chemical paper (also known as known print)
- GroundTruthedImage :The latent print is segmented manually giving each pixel a label of either shoeprint or background

Data Set Creation Process

• Latent print

- Person steps on powder and then on a carpet so that they leave a shoeprint on carpet
- Picture of the print is taken with a forensic scale near the print
- Current image resolution is calculated using the scale in the image and then scaled to 100dpi
- Chemical print
 - known print which is obtained by a person stamping on a chemical pad and then on a chemical paper, which would leave clear print on a paper.
 - Scanned into images of resolution 100dpi

Latent Print Database

- latent and chemical print pairs
- Latent images are cut into pieces to form partial images







Zappos Data Set

- Source: www.zappos.com
- Data set size : 10, 442 shoes
- Size on Disk : 1.08 GB
- Every Shoe has a side image, cross angle image, sole image and its name.





Timberland Urban Roll Top





Nike Tiempo Natural II TF





Converse All Star® Leather Hi









Converse All Star Leather Hi

Stacy Adams Taggart

New Balance MF825MK





Hi-Tec V-Lite Typhoon II



Bacco Bucci Hoover

Snipe Shoes 32436 (Mens)

PF Flyers Number 5

Rocky Mobilite

Computational Forensics National Academy of Sciences August 2008



$AI \rightarrow PR \rightarrow ML$

- From toy problems to real world
- From heuristics to statistically based machine learning systems
- Combination of computer science and statistical thinking

Science of PRML: Handwritten Digit Recognition



Wide variability of same numeral

- Handcrafted rules will result in large no of rules and exceptions
- Better to have a machine that learns from a large training set
- Provide a probabilistic prediction

Modern ML

- ML is a systematic approach
 - instead of ad-hoc methods in designing information processing systems
- There is much progress in ML theory
 - Unified theory: neural networks, SVM, Bayesian
 - New models, e.g., Gaussian processes
 - Fully Bayesian approach
 - Complex distributions, made feasible by sampling methods
 - Higher performing systems with small sample sizes
 - Course available on line
 - www.cedar.buffalo.edu/~srihari