# Role of Automation in the Forensic Examination of Handwritten Items

## Sargur Srihari
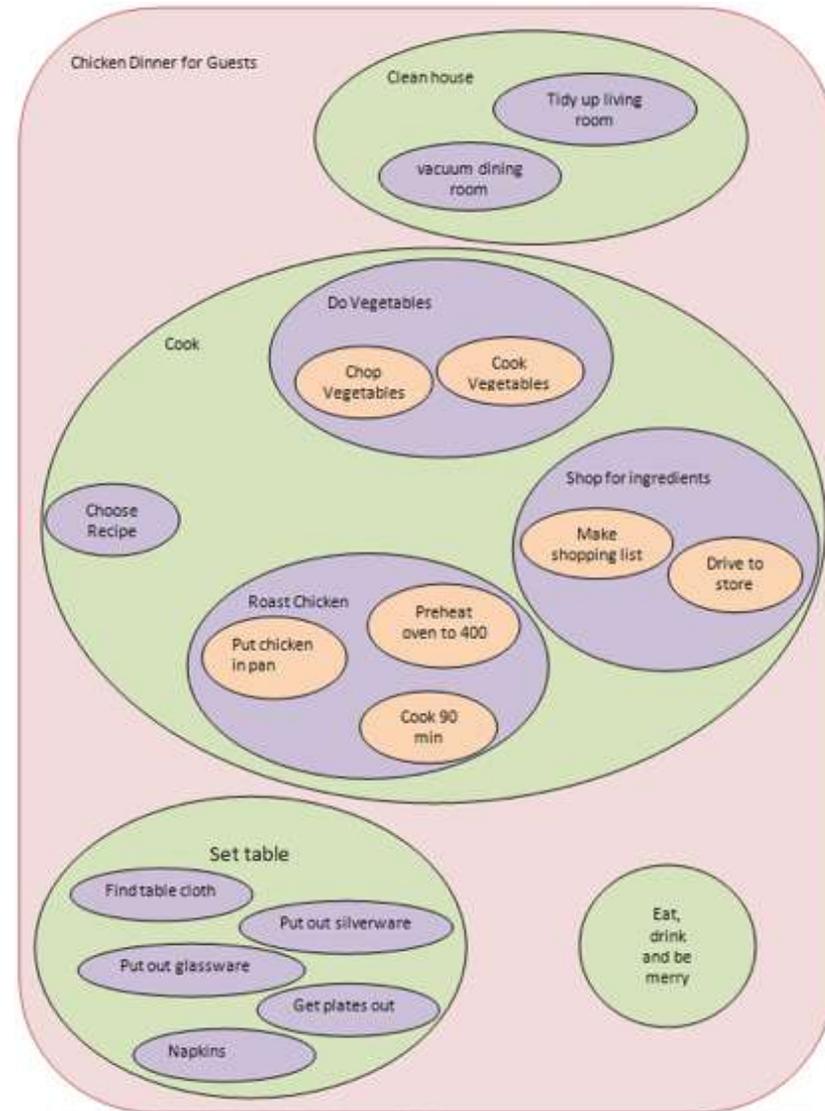
Department of Computer Science & Engineering

University at Buffalo, The State University of New York

NIST, Washington DC, June 3, 2013
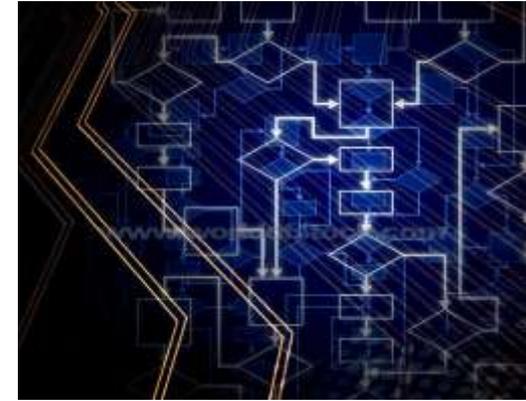
1

# Plan of Discussion

- Computational Thinking (CT)
  - What, Why, How, Limitations
- Reverse Engineering of QDE
- Automation Tools
  - Individualizing Characteristics
  - Opinion
  - Adequacy
- Summary and Discussion

# Computational Thinking (CT)

- ## What is it?

  - A way to solve problems, design systems, and understand human behavior

  - Draw on concepts of computer science

- ## Why?

  - To flourish in today's world, CT is the way to think and understand the world

S. Papert,   J. Wing

# How is CT done?

1. **Abstraction**
   – to understand and solve problems more effectively

2. **Algorithmic Thinking and Mathematics**
   – to develop efficient, fair, and secure solutions

3. **Understand scale**
   – Efficiency
   – Economic and social reasons

# CT and Law

- Long Dream: Logical rules to automate verdict
  - Napoleonic Code  (1804)
    - Minimize discretion, maximize predictability of outcome
    - Flounders: vagueness of words and variation of real world
  - Expert system replacements of judiciary
    - Poor record both of success and of uptake

- Better Inroads: Legal reasoning systems
  - Merely assist in legal decisions
    - E.g., Construct hypotheses for evidence in a crime scene
      - Remind detectives of hypotheses might have missed

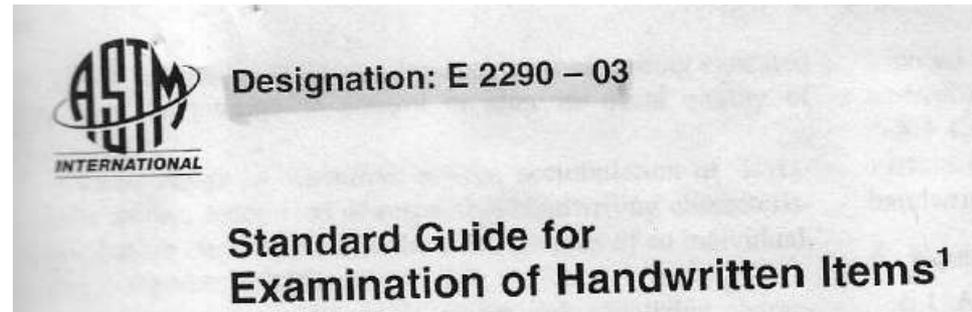- *Mind-expanding* avoids pitfalls of *mind-narrowing*

# CT and Forensics

- CT useful in domains where:
  - Human judgement is involved
  - Knowledge engineering can be performed
    - Starting point for creating artificial intelligence
- Within forensics:
  - Impression evidence
    - Handwriting, latent prints, footwear marks
- Handwriting:
  - Success demonstrated in recognition

# Knowledge Engineering for FDE

- CAT Principle
  - Comparability
  - Adequacy
  - Time Contemporaneous

Designation: E 2290 – 03

Standard Guide for
Examination of Handwritten Items[1]

- Characteristics
  - Class
  - Individualizing
  - Seven S's

Designation: E1658 – 08

Standard Terminology for
Expressing Conclusions of Forensic Document Examiners[1]

- Size, slant, spacing, shading, system, speed, strokes

# FDE- Exam. of HW Items (ASTM)

- Determine if  Q v Q,    K v K,    or     Q v K
- For Q and K:
  - Quality (copies?)
  - Distorted (disguised)
  - Type, Range
  - Individualizing characteristics?
- Comparable? else new K & repeat
- Differences/similarities for conclusion (5 or 9-pt)
  - Identification, Highly probable same, Probably did, Indications did , No conclusion Indications didn't, Probably didn't, Highly probable didn't, Elimination

# Word Cloud of ASTM Procedure

# Pseudo-code for Interactive Forensic Examination (iFOX)

---

**Algorithm 6** Comparison of handwritten items with statistical tools

---

1: *Determine Comparison Type*:
2:     $Q$ v $Q$ (no suspect or determine no. of writers)
3:     $K$ v $K$ (to determine variation range)
4:     $K$ v $Q$ (to determine/repudiate writership)
5: **for** each Q or K **do**
6:     *Quality:* determine visually or by automatic detection of noise.
7:     *Distortion:* detect manually or by use distortion measures.
8:     *Type determination:* manually or by automatic classification.
9:     *Internal consistency:* within document, e.g., multiple writers.
10:     *Determine range of variation*: compare subgroups.
11:     *Identify individualizing characteristics*: those with low probability.
12: **end for**
13: **for** each Comparison **do**
14:     *Comparability*: Both of same Type (Step 8).
15:     *Comparison*: Determine likelihood ratio (LR) based on characteristics and adequacy.
16:     *Form Opinion:* by quantizing LR or probability of identification.
17: **end for**

---

# Tools for Steps in FDE Procedure

- Quality

- Distortion

- Range

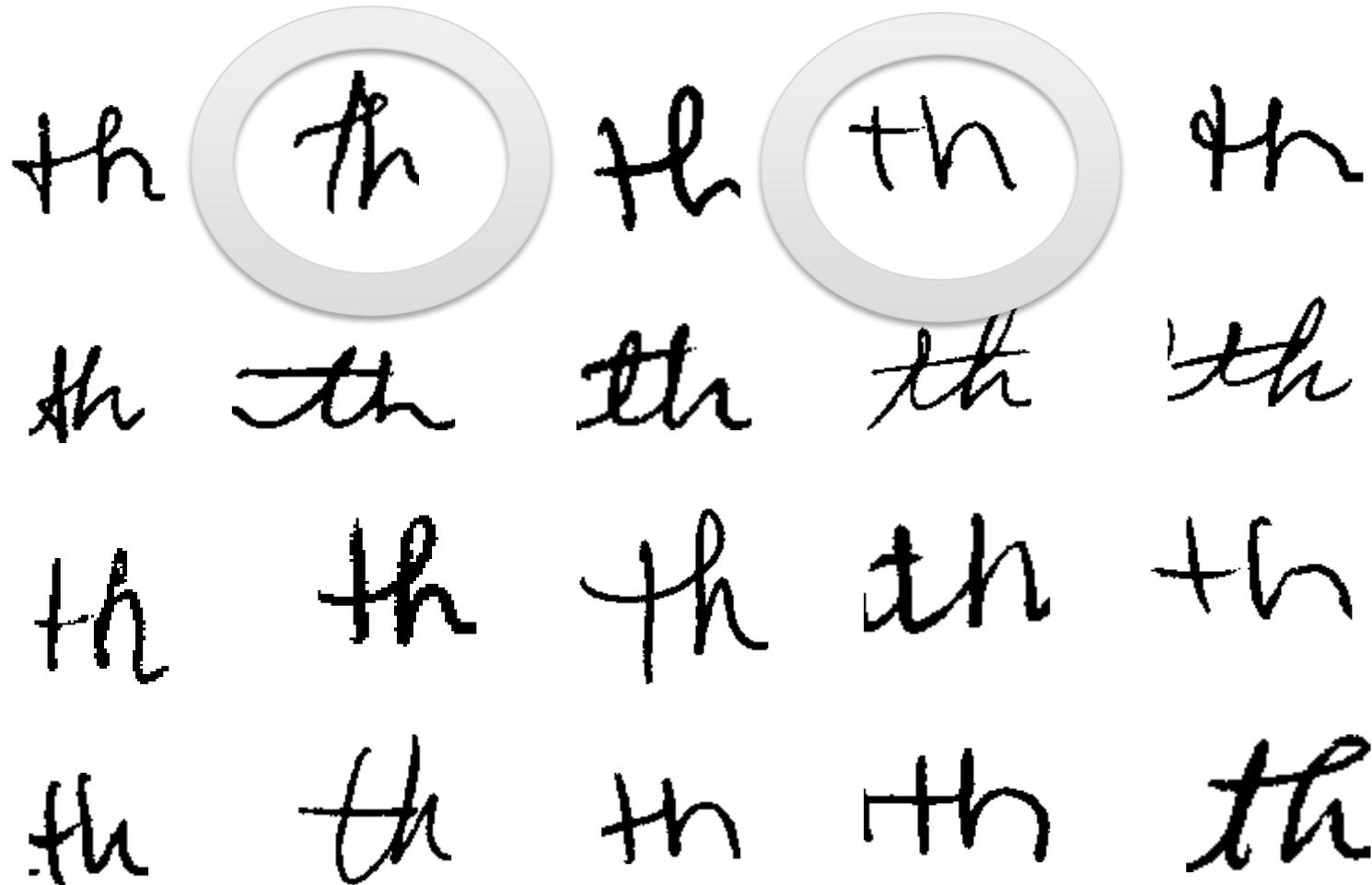1. Individualizing characteristics
2. Comparability (Type)
3. Comparison (Opinion)
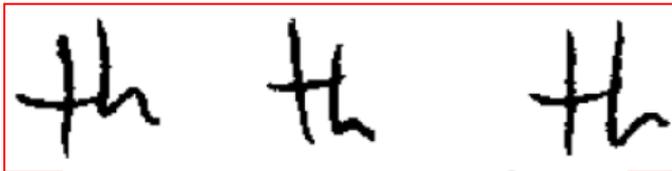4. Adequacy

Details
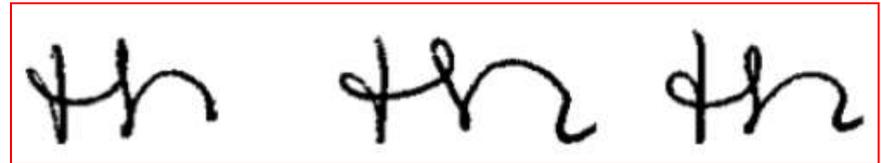Next

# Individualizing Characteristics?

# Characteristics of ``th''

| $R$ = **Height Relationship of $t$ to $h$** | $L$ = **Shape of Loop of $h$** | $A$ = **Shape of Arch of $h$** | $C$ = **Height of Cross on $t$ staff** | $B$ = **Baseline of $h$** | $S$ = **Shape of $t$** |
|---|---|---|---|---|---|
| $r^0$ = $t$ shorter than $h$ | $l^0$ = retraced | $a^0$ = rounded arch | $c^0$ = upper half of staff | $b^0$ = slanting upward | $s^0$ = tented |
| $r^1$ = $t$ even with $h$ | $l^1$ = curved right side and straight left side | $a^1$ = pointed | $c^1$ = lower half of staff | $b^1$ = slanting downward | $s^1$ = single stroke |
| $r^2$ = $t$ taller than $h$ | $l^2$ = curved left side and straight right side | $a^2$ = no set pattern | $c^2$ = above staff | $b^2$ = baseline even | $s^2$ = looped |
| $r^3$ = no set pattern | $l^3$ = both sides curved | | $c^3$ = no fixed pattern | $b^3$ = no set pattern | $s^3$ = closed |
| | $l^4$ = no fixed pattern | | | | $s^4$ = mixture of shapes |

R. J. Muehlberger, K. W. Newman, J. Regent and J. G. Wichmann, A Statistical Examination of Selected Handwriting Characteristics, *Journal of Forensic Sciences*, 1977: 206-210.
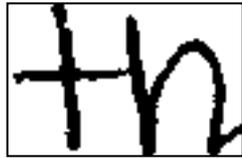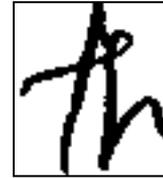


$$r^1, l^0, a^0, c^3, b^1, s^2$$



$$r^2, l^2, a^0, c^1, b^0, s^2$$
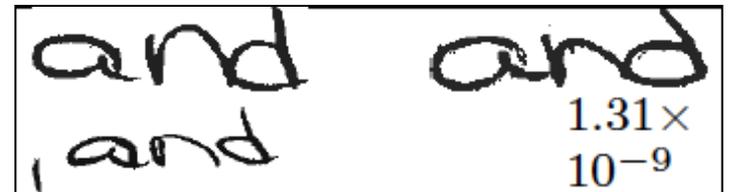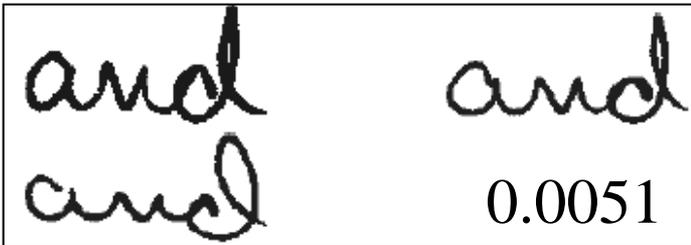
13

# Rarity: measure of Individualization

 0.0304

 $7.2 \times 10^{-8}$

High Probability
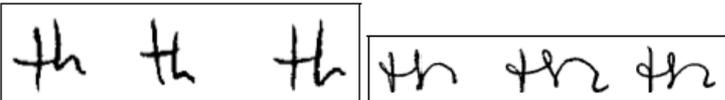
Low Probability

 0.0051

 $2.50 \times 10^{-9}$

 0.00167

 $1.31 \times 10^{-9}$

# Statistical Models of Characteristics



| R = Height Relationship of $t$ to $h$ | L = Shape of Loop of $h$ | A = Shape of Arch of $h$ | C = Height of Cross on $t$ staff | B = Baseline of $h$ | S = Shape of $t$ |
|---|---|---|---|---|---|
| $r^0$ = $t$ shorter than $h$ | $l^0$ = retraced | $a^0$ = rounded arch | $c^0$ = upper half of staff | $b^0$ = slanting upward | $s^0$ = tented |
| $r^1$ = $t$ even with $h$ | $l^1$ = curved right side and straight left side | $a^1$ = pointed | $c^1$ = lower half of staff | $b^1$ = slanting downward | $s^1$ = single stroke |
| $r^2$ = $t$ taller than $h$ | $l^2$ = curved left side and straight right side | $a^2$ = no set pattern | $c^2$ = above staff | $b^2$ = baseline even | $s^2$ = looped |
| $r^3$ = no set pattern | $l^3$ = both sides curved | | $c^3$ = no fixed pattern | $b^3$ = no set pattern | $s^3$ = closed |
| | $l^4$ = no fixed pattern | | | | $s^4$ = mixture of shapes |

# Probabilities for full joint distribution= 4,799

No. of Parameters if we assume independence= 19



| ...itial stroke of ...ormation of $a$ ($x_1$) | Formation of staff of $a$ ($x_2$) | Number of arches of $n$($x_3$) | Shape of arches of $n$ ($x_4$) | Location of mid-point of $n$($x_5$) | Formation of staff of $d$ ($x_6$) | Formation of initial stroke of $d$ ($x_7$) | Formation of terminal stroke of $d$ ($x_8$) | Symbol in place of t... word and ($x_9$) |
|---|---|---|---|---|---|---|---|---|
| ...ight of ...aff (0) | Tented (0) | One (0) | Pointed (0) | Above baseline (0) | Tented (0) | Overhand (0) | Curved up (0) | Formati...(0) |
| ...eft of ...aff (1) | Retraced (1) | Two (1) | Rounded (1) | Below baseline (1) | Retraced (1) | Underhand (1) | Straight across (1) | Symbol (1) |
| ...enter of ...aff (2) | Looped (2) | No fixed pattern (2) | Retraced (2) | At baseline (2) | Looped (2) | Straight across (2) | Curved down (2) | None (2 |
| ...o fixed ...attern ...) | No staff (3) | | Combination (3) | No fixed pattern (3) | No fixed pattern (3) | No fixed pattern (3) | No obvious ending stroke (3) | |
| | No fixed pattern (4) | | No fixed pattern (4) | | | | No fixed pattern (4) | |

# Probabilities=
287,999 (cursive)
809,999 (hand-print)

15

# What if we assume independence?

True Joint Probabilities:     Prob (height,weight)

| P(a,b) | b⁰ (heavy) | b¹ (light) | P(a) (height) |
|---|---|---|---|
| a⁰ (tall) | 0.6 | 0.05 | 0.65 |
| a¹ (short) | 0.05 | 0.3 | 0.35 |
| P(b) (weight) | 0.65 | 0.35 | |

P(tall, light) < P(short,light)    0.05<0.3
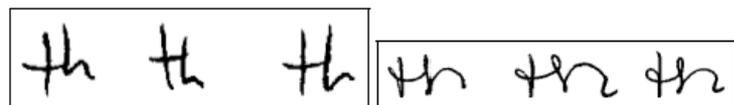Short & light six times more probable  than tall and light: Correct!

## Assuming Independence

| P(a,b) | b⁰ (heavy) | b¹ (light) | P(a) (height) |
|---|---|---|---|
| a⁰ (tall) | 0.42 | 0.23 | 0.65 |
| a¹ (short) | 0.23 | 0.12 | 0.35 |
| P(b) (weight) | 0.65 | 0.35 | |

P(tall,light) > P(short,light)  0.23 >0.12
Tall & light, twice probability of short & light: Wrong!
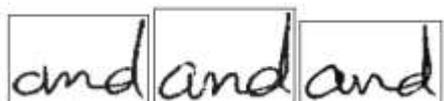
# Compromise Solution: PGMs

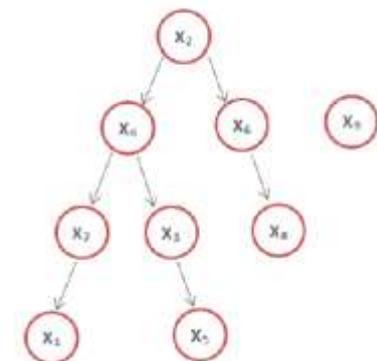| R = Height Relationship of $t$ to $h$ | L = Shape of Loop of $h$ | A = Shape of Arch of $h$ | C = Height of Cross on $t$ staff | B = Baseline of $h$ | S = Shape of $t$ |
|---|---|---|---|---|---|
| $r^0$ = $t$ shorter than $h$ | $l^0$ = retraced | $a^0$ = rounded arch | $c^0$ = upper half of staff | $b^0$ = slanting upward | $s^0$ = tented |
| $r^1$ = $t$ even with $h$ | $l^1$ = curved right side and straight left side | $a^1$ = pointed | $c^1$ = lower half of staff | $b^1$ = slanting downward | $s^1$ = single stroke |
| $r^2$ = $t$ taller than $h$ | $l^2$ = curved left side and straight right side | $a^2$ = no set pattern | $c^2$ = above staff | $b^2$ = baseline even | $s^2$ = looped |
| $r^3$ = no set pattern | $l^3$ = both sides curved | | $c^3$ = no fixed pattern | $b^3$ = no set pattern | $s^3$ = closed |
| | $l^4$ = no fixed pattern | | | | $s^4$ = mixture of shapes |

$$P(X) = P(R)P(L|S)P(A|L)P(C|S)P(B|R,A)P(S|R)$$

100 parameters

| Initial stroke of formation of $a$ ($x_1$) | Formation of staff of $a$ ($x_2$) | Number of arches of $n$ ($x_3$) | Shape of arches of $n$ ($x_4$) | Location of midpoint of $n$ ($x_5$) | Formation of staff of $d$ ($x_6$) | Formation of initial stroke of $d$ ($x_7$) | Formation of terminal stroke of $d$ ($x_8$) | Symbol in place of the word *and* ($x_9$) |
|---|---|---|---|---|---|---|---|---|
| Right of staff (0) | Tented (0) | One (0) | Pointed (0) | Above baseline (0) | Tented (0) | Overhand (0) | Curved up (0) | Formation (0) |
| Left of staff (1) | Retraced (1) | Two (1) | Rounded (1) | Below baseline (1) | Retraced (1) | Underhand (1) | Straight accross (1) | Symbol (1) |
| Center of staff (2) | Looped (2) | No fixed pattern (2) | Retraced (2) | At baseline (2) | Looped (2) | Straight accross (2) | Curved down (2) | None (2) |
| No fixed pattern (3) | No staff (3) | | Combination (3) | No fixed pattern (3) | No fixed pattern (3) | No fixed pattern (3) | No obvious ending stroke (3) | |
| | No fixed pattern (4) | | No fixed pattern (4) | | | | No fixed pattern (4) | |

99 parameters (cursive)
77 parameters (hand-print)
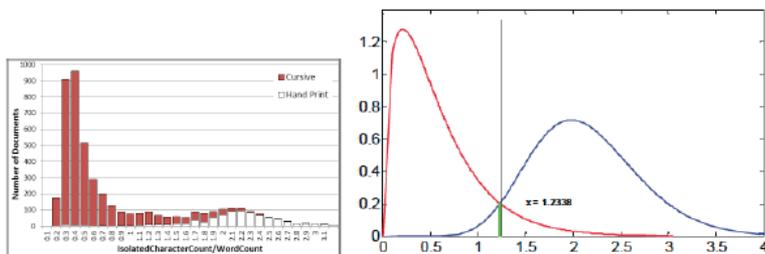
# Learning PGMs from Data

- Bayesian Networks (directed graphs)
- Markov Networks (undirected graphs)
- Learning algorithms:
  - Determine pairwise independences using chi-squared tests
  - Determine quality of model using log-loss
  - Problem is NP-hard
    - use approximate solutions

# Type Determination

- ## Cursive vs. Handprint



$f_1$: Discreteness
   Ratio of isolated character count(ICC)  to word count (WC)
$f_2$: Loopiness
   Ratio of interior to exterior contours

# Opinion

$$LR_J = LR(\mathbf{k}, \mathbf{q}) = \frac{P(\mathbf{k}, \mathbf{q}|h^0)}{P(\mathbf{k}, \mathbf{q}|h^1)}$$

$$LR_D = \frac{P(D(\mathbf{k}, \mathbf{q})|h^0)}{P(D(\mathbf{k}, \mathbf{q})|h^1)}$$
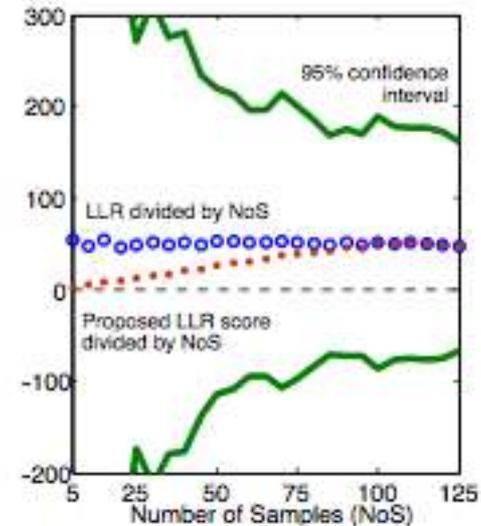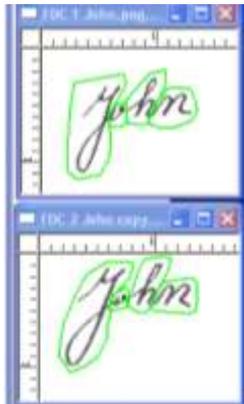
$$P(h^0|F) = \frac{P(h^0) \cdot \prod_i LR(f_i)}{P(h^1) + P(h^0) \cdot \prod_i LR(f_i)}$$

| Information $\mathcal{I}$ | Content $\mathcal{C}$ | Sys. Accuracy | Min.LLR | Max.LLR |
|---|---|---|---|---|
| Line | Same | 86.40% | -93.82 | 115.57 |
| | Different | 62.98% | -72.14 | 11.05 |
| Multiple lines | Same | 93.81% | -105.02 | 96.83 |
| | Different | | | |
| Half page | Same | 93.08% | -322.59 | 698.64 |
| | Different | 94.78% | -111.83 | 172.28 |
| Full page | Same | 95.75% | -90.1 | 67.93 |

| Scale | Opinions for same | $P_{ic}^S$ |
|---|---|---|
| 1 | Identified as same | $0.00 \sim 22.21$ |
| 2 | Highly probably same | $22.22 \sim 44.43$ |
| 3 | Probably same | $44.44 \sim 66.65$ |
| 4 | Indicating same | $66.66 \sim 88.87$ |
| 5 | No conclusion | $88.88 \sim 100.00$ |

| Scale | Opinions for different | $P_{ic}^{Df}$ |
|---|---|---|
| 5 | No conclusion | $88.88 \sim 100.00$ |
| 6 | Indicating different | $66.66 \sim 88.87$ |
| 7 | Probably different | $44.44 \sim 66.65$ |
| 8 | Highly probable different | $22.22 \sim 44.43$ |
| 9 | Identified as different | $0.00 \sim 22.21$ |

20

# Adequacy





1. A single feature $F = f_1$ with $LR(f_1) = 96$.

2. Nine features $F = \{f_i\}_{i=1}^{9}$ with $\{LR(f_i), i = 1, .., 9\} = \{3, 4, 2, \frac{1}{4}, 2, 2, \frac{1}{3}, 6, 2\}$
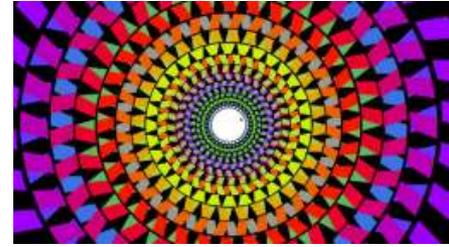
# Availability of Automation Tools

- Interactive tools rather than automation
- Incorporate CT
  - Abstraction: Aids to organize thought process
  - Algorithms and Mathematics
  - Scalability
    - Potential of Large quantities quickly analyzed
- Mind Expanding
  - Probability allows considering characteristics otherwise ignored, or discounting identified ones
  - Value of small amounts of information

# Status of Automation Tools

- Interactive tools rather than automation
- Work in Progress
  - Characteristics
    - Data needs to be collected
  - Learning statistical models
    - Learning PGMs is current topic in ML
  - Inference algorithms
  - Type determination
  - Opinion Mapping

23

# Summary

- Computational Thinking + Forensics = Computational Forensics
- Solve using
  - Abstraction
  - Algorithms
  - Mathematics
  - Scale

# Summary

- Reverse engineering of QDE
  - Available in ASTM standards, other QD literature
- Steps amenable to automation tools
  - Data Collection
  - Modeling distribution of characteristics
  - Type determination
  - Likelihood Ratios (Opinion)
  - Confidence Intervals (Adequacy)

# How does this fit with fully Automated  Systems?

- Systems such as FISH, CEDAR-FOX and FLASH-ID narrow down possibilities

- Interactive systems (iFOX) will assist the document examiner in going the last mile
  - E.g., associate  probabilities with their observations