

A Statistical Approach to Line Segmentation in Handwritten Documents

M. Arivazhagan, H. Srinivasan, S.N.Srihari
Center of Excellence for Document Analysis
and Recognition (CEDAR)
University at Buffalo
State University of New York

Overview

1. Algorithm Objective
2. Algorithm Steps
 1. Initial candidate lines – Projection profile
 2. Statistical Modeling and Decision Making
 3. Line segmentation
3. Results
 1. Different Data Sets
 2. Comparison with Previous Algorithm
4. Conclusion

Algorithm Objective

Simple image

Highly skewed image

Complex image
with overlapping
lines

From
Jim Elder
829 Loop Street, Apt 300
Alhambra, New York 14707

To
Dr. Bob Grant
602 Queensbury Parkway
Omni, West Virginia 25638

Nov 10, 1999

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

It all started around 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, xrays and blood tests later, we were told it was just exhaustion.

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

Thank you!
Jim

Nov 10, 1999

From
Jim Elder
829 Loop Street, Apt 300
Alhambra, New York 14707

To
Dr. Bob Grant
602 Queensbury Parkway
Omni, 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 around 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, xrays and blood tests later, we were told it was just exhaustion.

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

Thank you!
Jim

Nov 10, 1999

From
Jim Elder
829 Loop Street, Apt 300
Alhambra, New York 14707

To
Dr. Bob Grant
602 Queensbury Parkway
Omni, 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 around 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, xrays and blood tests later, we were told it was just exhaustion.

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

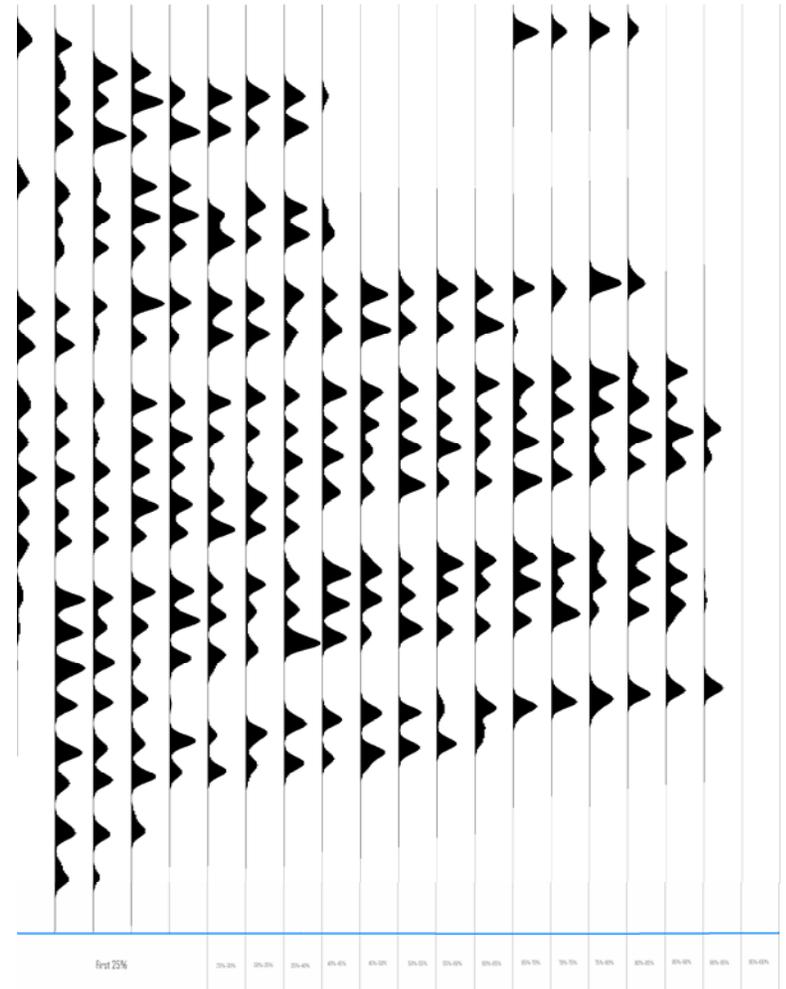
Thank you!
Jim

Algorithm Outline

- Thresholding and Chain code Representation
- Initial candidate lines
 - Local Projection Profiles
- Line drawing algorithm
 - Bi-variate probability densities to model line
 - Probabilistic decision to classify
 - Lines guided by piece wise projection profile

Local Projection Profiles

- Smoothed projection profile for each 5% chunk.
- Smoothed projection profile of first 25% of document used as initial lines.
- Remaining projection profiles help guide the lines.



Initial candidate lines

Nov 10, 1979

From
Jim Elder
621 Loop Street Apt 311
Albany, New York 12207

To
Dr. Bob Grant
602 Queensberry Parkway
Charlottesville, West Virginia 25408

We were referred to you by Vera Cohen at the University Medical Center. This is regarding my ~~friend~~ ^{friend} Kate Zick.

It all started ^{about} ~~about~~ six months ago while attending the "Yale" Jazz Concert. Organizing such an event is no ~~thing~~ ^{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.

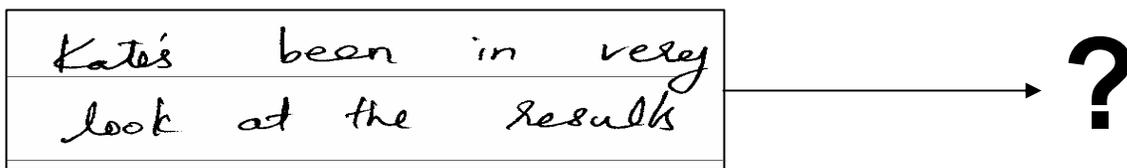
However, the extra hours affected her health; halfway through the show she passed out. We rushed her to the hospital, and several questions I say and blood tests later, were told it was just exhaustion.

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

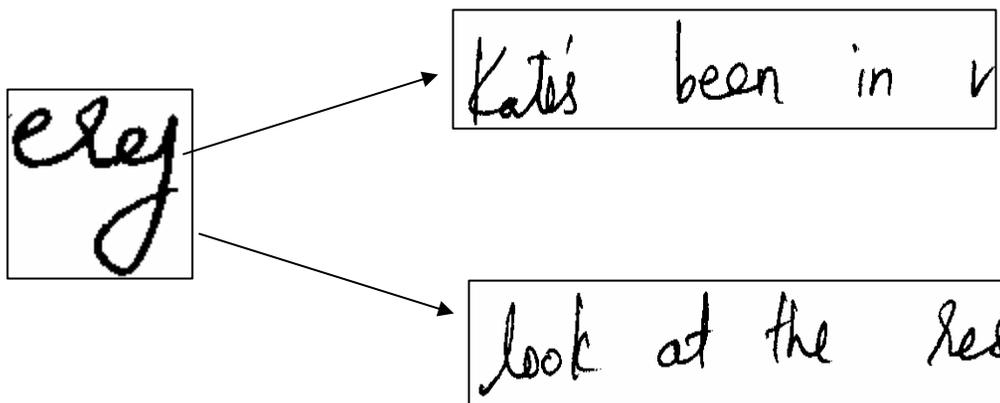
Thank you!

Jim

Ambiguous Component Decision



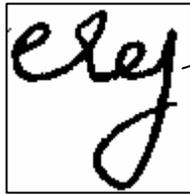
COMP



Distributions of Line Structures

STEP 1

COMP
'N' points



Kate's been in v

Line_Above $\sim N(\mu_a, \Sigma_a)$

$\mu_a = \{404.13, 2548.97\}$

$\Sigma_a = \begin{bmatrix} 28944.46 & -546.58 \\ -546.58 & 208.59 \end{bmatrix}$

look at the res

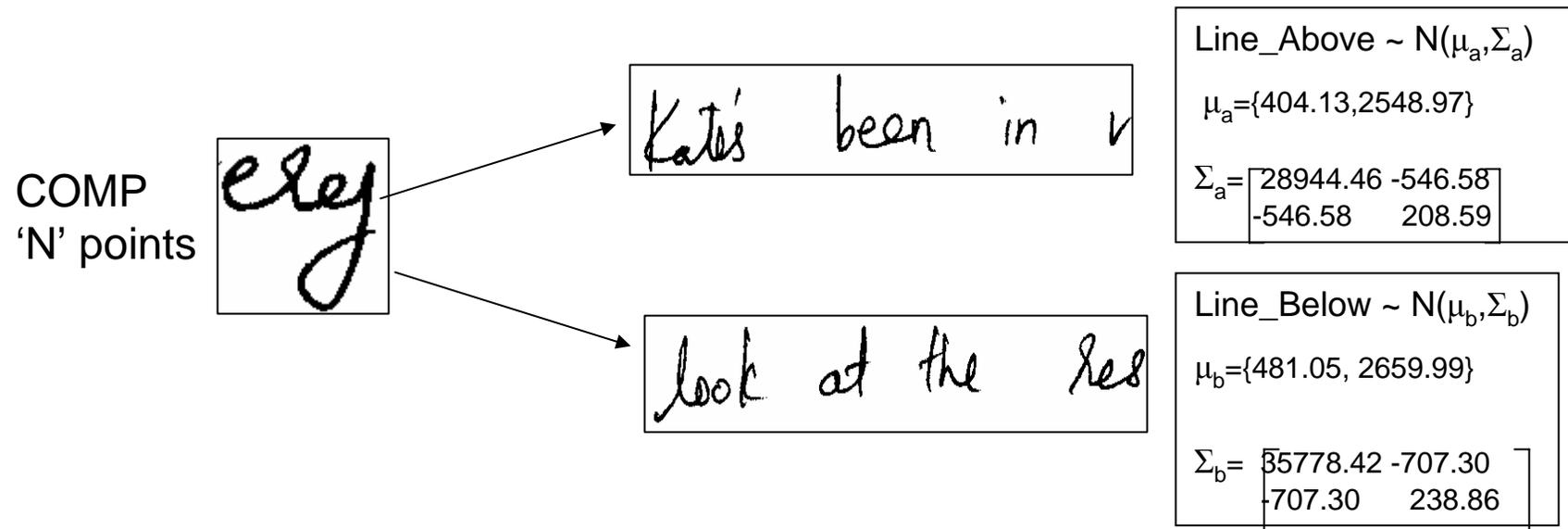
Line_Below $\sim N(\mu_b, \Sigma_b)$

$\mu_b = \{481.05, 2659.99\}$

$\Sigma_b = \begin{bmatrix} 35778.42 & -707.30 \\ -707.30 & 238.86 \end{bmatrix}$

Line Probabilities

STEP 2



$$P(C | \vec{\mu}_A, \Sigma_A) = P(p_1 | \vec{\mu}_A, \Sigma_A) \cdot P(p_2 | \vec{\mu}_A, \Sigma_A, p_1) \cdots P(p_T | \vec{\mu}_A, \Sigma_A, p_1, p_2, \dots, p_{T-1})$$

$$\log P(C | \text{Line_Above}) = -22161.54$$

$$\log P(C | \text{Line_Below}) = -44391.26$$

Likelihood Ratio Decision

STEP 1

COMP
'N' points

very

Kate's been in v

look at the res

Line_Above $\sim N(\mu_a, \Sigma_a)$

$\mu_a = \{404.13, 2548.97\}$

$\Sigma_a = \begin{bmatrix} 28944.46 & -546.58 \\ -546.58 & 208.59 \end{bmatrix}$

Line_Below $\sim N(\mu_b, \Sigma_b)$

$\mu_b = \{481.05, 2659.99\}$

$\Sigma_b = \begin{bmatrix} 35778.42 & -707.30 \\ -707.30 & 238.86 \end{bmatrix}$

$$P(C | \vec{\mu}_A, \Sigma_A) = P(p_1 | \vec{\mu}_A, \Sigma_A) \cdot P(p_2 | \vec{\mu}_A, \Sigma_A, p_1) \dots P(p_T | \vec{\mu}_A, \Sigma_A, p_1, p_2, \dots, p_{T-1})$$

STEP 2 $\log P(C | \text{Line_Above}) = -22161.54$

$\log P(C | \text{Line_Below}) = -44391.26$

STEP 3

Kate's been in very
look at the results

Distance metric decision

- When component in first 25% of document
 - Very few pixels to model line statistics
- No foreground pixels exist in the above/below line
- Strength of line evidence of Gaussian is poor.

$$p_u = \frac{|y_h - y_u|}{|y_h - y_u| + |y_h - y_d|} \quad \text{and} \quad p_d = \frac{|y_h - y_d|}{|y_h - y_u| + |y_h - y_d|}.$$

$p_h = [x_h, y_h]$: Pixel at which the line hits the component
 y_u : Pixel of minimum y co-ordinate in the upward traversal.
 y_d : Pixel of maximum y co-ordinate in the downward traversal.

If p_u is greater, component belongs to above line, else below line

Overlapping lines

Nov 10, 1979

From
Jim Eldor
829 Loop Street, Apt 300
Allentown, New York 14707

To
Dr. Bob Grant
602 ~~Queensberry Parkway~~
Omar, ~~West Virginia~~ 25638

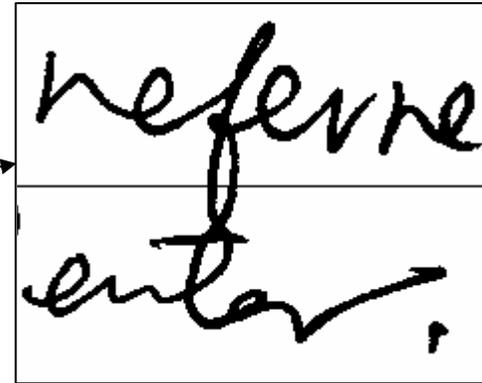
We were referred to you by Xeno Chorn at the University Medical Center. This is regarding my friend, Kate Zack.

It all started around six months ago while attending the "Kubies" Jazz Concert. ~~Organizing~~ such an event was no picnic, and as President of ~~the Alumni Association~~, a co-sponsor of the event, Kate was overwhelmed. 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.

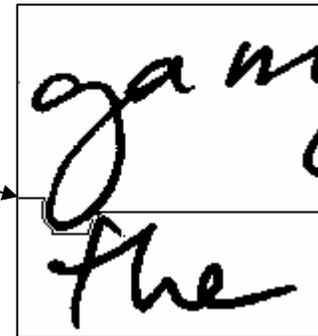
Kate's been in very bad health since. Could you kindly take a look at the results and give us your opinion?

Thank you!
Jim



referred
enter,

Cut through overlapping component



refer
the
enter

Traversal

Ongoing research addresses
this problem

refer
enter

refer

enter

Experimental Results:

Comparison with data sets, algorithms

Database	# of lines	Previous method's (clustering based) accuracy	Proposed method's accuracy
CEDAR-LETTER 300 images	7201	88.86 %	97.01 %
Arabic 120 images	2030	87.33 %	98.62 %
School exam essay 300 images	2350	86.43 %	96.30 %

Experimental Results: Line and Component Level Accuracy

	Total	Number correct	Accuracy
Number of lines	4768	4656	97.65 %
Number of Components	78902	77958	98.81 %
Cut-through	928	680	74.36 %

Summary

- New approach to handwritten line segmentation
 - Global properties (projection profile)
 - Statistics of line structure
 - Local rules when not enough data
- Algorithm robust to handle skewed documents
- Accurate association of components to correct lines – 99%
- Improved separation of overlapping contours still a research issue