# Keyword Searching and Indexing in Large Collections of Handwritten Documents

**Roger Labahn**

*Computational Intelligence Technology Lab*

Mathematical Optimization
Institute for Mathematics

University of Rostock

Germany

Prologue

Foundation

Application

Epilogue

## Question

Do we really need a good transcription for searching and investigating?

## Ambitions

- understand different concept & its applications
- learn technological terminology
- know about configuration & behaviour, features & bugs, …

## Note !

- tools for continuous work – NOT just engines for execution
- adapt for specific challenges – understand & interpret outcomes

Foundation

## Text Recognition Process

- Recognition Engine ⤳ *character confidence scores* per position
- Neural Network / HMM outputs: estimate character probabilities

## Raw Reading Result

- method: choose most likely character per position
- *free* reading: without considering document context
  e.g.: language, time, writer, …

## Post – OCR Correction

- find & correct errors: use external sources from document context
  e.g.: language models, dictionaries, transcripts, …
- ⤳ strings: text transcription

## Text Recognition Process

- Recognition Engine ⤳ *character confidence scores* per position
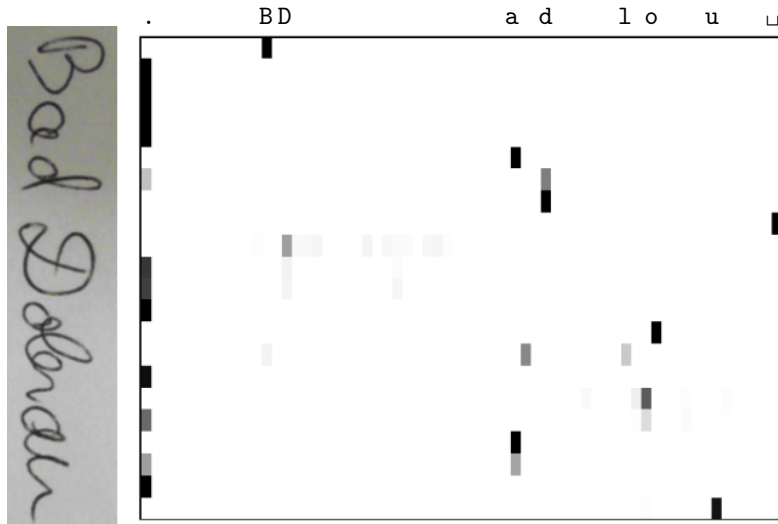- Neural Network / HMM outputs: estimate character probabilities

## Confidence Matrix – CONFMAT

- idea: evaluate entire recognition information
- application: Store this text recognition result!

## Decoding

- query strings: use external sources from document context
- find optimal match / representation: query $\Longleftrightarrow$ ConfMat
- ⤳ ranked alternatives

# Confidence Matrix

# Measuring Similarity: String vs. String

## Definition (LEVENSHTEIN Distance)

$$\mathrm{dist}(\mathrm{string1}, \mathrm{string2}) := \mathrm{count}(\mathrm{insertions}, \mathrm{deletions}, \mathrm{substitutions})$$

## Algorithm        Dynamic Programming

- extremely efficient – very fast
- finds optimal (shortest / cheapest) path through weight (distance / cost) matrix
- weights: distance ⟿ cost
- counting ⟿ adding costs (weights)

**Example:** dist(WIEN, WEIN) = 2

### insertions & deletions

|   |   | W | E | I | N |
|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 |
| W | 1 | 0 | 1 | 2 | 3 |
| I | 2 | 1 | 2 | 1 | 2 |
| E | 3 | 2 | 1 | 2 | 3 |
| N | 4 | 3 | 2 | 3 | 2 |

W   I E N
W E I   N

### with substitutions

|   |   | W | E | I | N |
|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 |
| W | 1 | 0 | 1 | 2 | 3 |
| I | 2 | 1 | 1 | 1 | 2 |
| E | 3 | 2 | 1 | 2 | 2 |
| N | 4 | 3 | 2 | 2 | 2 |

W I E N
W E I N

▸ Demo

# Measuring Similarity: String vs. ConfMat

| . | B | a | d |  | D | o | b | e | r | a | n |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 1.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 1.00 | 0.00 | 1.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 1.00 | 0.00 | 1.00 | 1.99 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 1.00 | 0.00 | 1.00 | 1.99 | 2.99 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 1.00 | 0.00 | 0.99 | 1.99 | 2.99 | 3.99 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 1.99 | 1.00 | 0.00 | 1.99 | 2.99 | 3.99 | 4.98 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 2.30 | 1.31 | 0.31 | 0.00 | 2.64 | 3.64 | 4.64 | 5.64 | 0.00 | 0.00 | 0.00 | 0.00 |
| 3.30 | 2.30 | 1.31 | 0.00 | 1.00 | 3.64 | 4.64 | 5.64 | 6.63 | 0.00 | 0.00 | 0.00 |
| 4.30 | 3.30 | 2.30 | 1.00 | 0.00 | 1.99 | 4.64 | 5.64 | 6.63 | 7.63 | 0.00 | 0.00 |
| 4.81 | 3.81 | 2.81 | 1.51 | 0.51 | 0.00 | 2.50 | 5.15 | 6.15 | 7.14 | 8.14 | 0.00 |
| 4.81 | 3.81 | 2.81 | 1.51 | 0.51 | 0.00 | 0.90 | 3.40 | 7.05 | 7.05 | 8.04 | 9.04 |
| 4.81 | 3.81 | 2.81 | 1.51 | 0.51 | 0.00 | 0.87 | 1.77 | 4.28 | 6.92 | 7.92 | 8.92 |
| 0.00 | 3.81 | 2.81 | 1.51 | 0.51 | 0.00 | 0.87 | 1.77 | 2.77 | 5.27 | 7.91 | 8.91 |
| 0.00 | 0.00 | 3.81 | 2.50 | 1.50 | 0.99 | 0.00 | 1.87 | 2.77 | 3.76 | 6.26 | 8.91 |
| 0.00 | 0.00 | 0.00 | 3.11 | 2.11 | 1.61 | 0.61 | 0.00 | 2.48 | 3.38 | 4.37 | 6.88 |
| 0.00 | 0.00 | 0.00 | 0.00 | 2.11 | 1.61 | 0.61 | 0.00 | 0.99 | 3.38 | 4.37 | 5.36 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 2.39 | 1.39 | 0.78 | 0.79 | 1.76 | 4.17 | 4.37 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.39 | 0.78 | 0.79 | 1.49 | 2.48 | 4.37 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.78 | 1.78 | 1.78 | 1.49 | 2.48 | 3.48 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.78 | 1.78 | 1.78 | 1.49 | 2.00 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.78 | 1.78 | 1.49 | 2.00 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 2.46 | 2.45 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 2.45 |

# Probability / Confidence ⤳ Distance / Cost

| Confidence | | Distance |
|---|---|---|
| 1.00 | $\exp(-\text{dist})$ | 0.0 |
| ⋮ | $\Longleftarrow$ | ⋮ |
| 0.00 | | ⋮ |
| | $\Longrightarrow$ | ⋮ |
| | $-\ln(\text{prob})$ | $\infty$ |

## Asymptotic Problem

- confidence mapping / scaling of arbitrarily large distances
- What distances practically correspond to probability Zero?

Prologue

Foundation

Application

Epilogue

## Setup

### Text Recognition Engine

Text $\rightsquigarrow$  ConfMats
per line

### Decoding Engine

Query $\rightsquigarrow$  Alternatives ranked by
confidence / distance

### Measurement Issues

- Confidences are NO PROBABILITIES!
- Confidences / Distances have NO ABSOLUTE meaning!
- Measurements are essentially INCOMPARABLE across different queries!
- Thresholds require MANUAL CONFIGURATION & TUNING!

## Keyword Search

### String Search

INPUT   Query String

Decoding   Everywhere: distance to query

OUTPUT   rank "reasonable" hits – skip "irrelevant" answers

### Big Data Issues

- Response time: inacceptable ! ?
- Preprocessing: index ⤳ database

▸ UPVLC

## Reading Text

INPUT Language Model – Dictionary incl. word frequencies

Decoding distances: all dictionary entries everywhere

OUTPUT "reasonable" text alternatives: close to ConfMat & Language Model

## Issue

- both poor Text Recognition AND Language Model

# Investigating

## Challenges

- demanding Language Models: fuzzy / dynamic / incomplete
- complex queries:
  specific combinations – character classes – restricted vocabulary

## Regular Expression Decoding

- *Regular Expression*: Computer Science & Programming
- unsupported features: named classes, …
- additional feature: dictionary classes

# Investigation with Regular Expressions

## Regular Expression Example

four-digit year:     1YYY

1[0-9][0-9][0-9]

1[0-9]{3}

## Regular Expression Syntax

### Regular Expression Example

complete date:    `TT.MM.YYYY`

`.*(?<KW>[0-3][0-9]\.[0-1][0-9]\.[1-2][0-9]{3}).*`

### Note

- match: against the entire line
- score: designated KW-group

`.*(?<KW><query>).*`

# Performance

## TRANSKRIBUS KWS Expert Mode

- Regular Expression Decoding: `.*(?<KW>[0-9]{4}).*`
- KWS demo collection: 388 pages from StAZh 1796-97
- » $\approx 40\,s$

## Note

- strongly depends upon hardware
- coming up next:
  massive parallelization on CPU & GPU

## Performance

**Office laptop    single core**

- Searching 1 keyword in 10.500 lines (433 BENTHAM pages):
- » 2 – 3 s average

- Reading 1 page against 11.650 words dictionary:
- » 8 – 9 s average

## Questions

### Your wishes – expectations – requirements . . .

- . . . realistic query type / elaborateness / complexity ??

- . . . realistic data corpus size ??

- . . . realistic query response time ??

# THANKS …

## CIT lab   Group

## MoU   Partner   SME

PLANET artificial intelligence GmbH

## EU   HORIZON2020   Grant

**READ** Recognition and Enrichment of Archival Documents

… for your kind attention!