

The Long and the Short of It

summarising event sequences with serial episodes

Nikolaj Tatti & **Jilles Vreeken**



Question of the day

How can we discover
the key patterns from an event sequence?



Summarising Event Sequences

The **ideal** outcome of pattern mining

- patterns that show the structure of the data
- preferably a small set, without redundancy or noise

Frequent pattern mining does **not** achieve this

- pattern explosion → overly many, overly redundant results


We pursue the ideal for serial episodes

- we want a group of patterns that summarise the data well
- we take a **pattern set** mining approach

Event sequences

Alphabet Ω $\{ a, b, c, d, \dots \}$

Data D $a \ b \ d \ c \ a \ d \ b \ a \ a \ b \ c \ a \ d \ a \ b \ a \ b \ c$

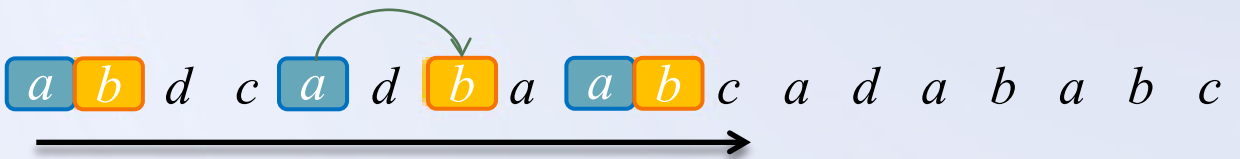


one, or
multiple
sequences

$\{ a \ b \ d \ c \ a \ d \ b \ a \ a \ b \ c ,$
 $a \ b \ d \ c \ a \ d \ b ,$
 $a \ b \ d \ c \ a \ d \ b \ a \ a , \dots \}$

Event sequences

Alphabet Ω $\{ a, b, c, d, \dots \}$

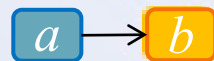
Data D 

one, or
multiple
sequences

$\{ a b d c a d b a a b c,$
 $a b d c a d b,$
 $a b d c a d b a a, \dots \}$

Patterns

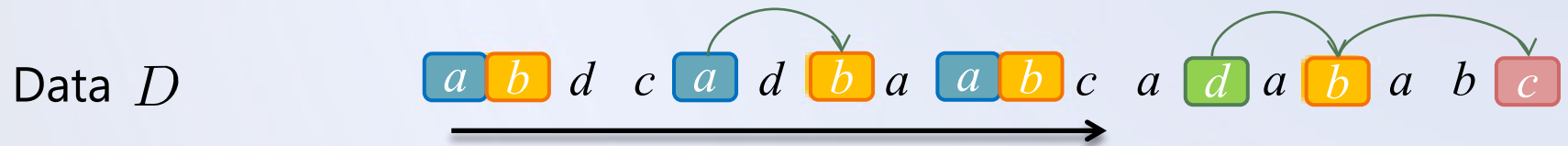
serial
episodes



'subsequences
allowing gaps'

Event sequences

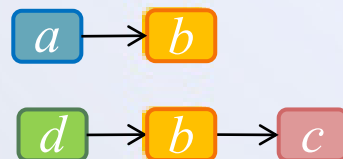
Alphabet Ω $\{ a, b, c, d, \dots \}$



$\{ a b d c a d b a a b c, \\ a b d c a d b, \\ a b d c a d b a a, \dots \}$

Patterns

serial episodes



'subsequences allowing gaps'

Summarising Event Sequences

We want to find good summaries.

Three important questions

1. how do we **score** a pattern-based summary?
2. how do we **describe** a sequence given a pattern set?
3. how do we **find** good pattern sets?

Summarising Event Sequences

We want to find good summaries.

Three important questions

1. how do we **score** a pattern-based summary?
2. how do we **describe** a sequence given a pattern set?
3. how do we **find** good pattern sets?

Scoring a Summary

We want models that generalise the data

and hence, we need a score that

- **rewards** models that identify real structure, and
- **punishes** redundancy and noise

No off-the-shelf score available for serial episodes

- e.g. no well-founded priors
- we can, however, make these goals concrete by **MDL**

MDL

The Minimum Description Length (MDL) principle

given a set of models \mathcal{M} , the best model $M \in \mathcal{M}$
is that M that minimises

$$L(M) + L(D|M)$$

in which

$L(M)$ is the length, in bits, of the description of M

$L(D|M)$ is the length, in bits, of the description of
the data when encoded using M

MDL for Event Sequences

By MDL we define

the optimal set of serial episodes as the set that describes the data most succinctly

To use MDL, we need

- a lossless encoding for our models,
- a lossless encoding for the data given a model


Models

As models we use **code tables**

- dictionaries of patterns and associated codes





| <i>pattern</i> | <i>code</i> | <i>gap</i> | <i>non-gap</i> |
|----------------|-------------|------------|----------------|
| <i>abc</i> | <i>p</i> | <i>?</i> | <i>!</i> |
| <i>da</i> | <i>q</i> | <i>?</i> | <i>!</i> |
| <i>a</i> | <i>a</i> | | |
| <i>b</i> | <i>b</i> | | |
| <i>c</i> | <i>c</i> | | |
| <i>d</i> | <i>d</i> | | |

Encoding Event Sequences

Data D : $a \ b \ d \ c \ a \ d \ b \ a \ a \ b \ c$


Encoding 1: using only singletons

C_p 

CT_1 : a 
 b 
 c 
 d 

The length of the code \boxed{X} for pattern X

$$L(\boxed{X}) = -\log(p(\boxed{X})) = -\log\left(\frac{usg(X)}{\sum usg(Y)}\right)$$

The length of the code stream

$$L(C_p) = \sum_{X \in CT} usg(X) L(\boxed{X})$$

Encoding Event Sequences

Data D : $a \ b \ d \ c \ a \ d \ b \ a \ a \ b \ c$

Encoding 2: using patterns

C_p $p \ d \ a \ q \ b \ p$

C_g $! \ ? \ ! \ ? \ ! \ ! \ !$

gap

gap

Alignment: $a \ b \ d \ c \ a \ d \ b \ a \ a \ b \ c$

$p \ ! \ ? \ !$ $q \ ? \ !$ $p \ ! \ !$

CT_2 :


| | | | |
|-------|-----|-----|-----|
| a | a | | |
| b | b | | |
| c | c | | |
| d | d | | |
| abc | p | $?$ | $!$ |
| da | q | $?$ | $!$ |

gaps

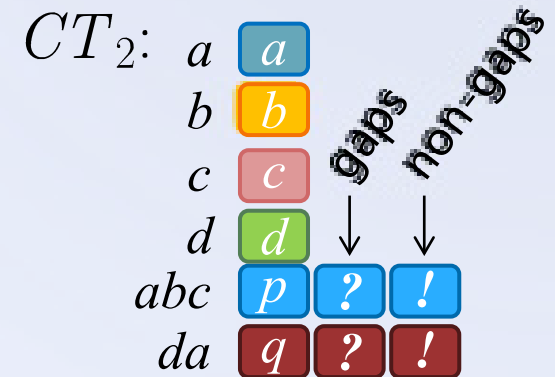
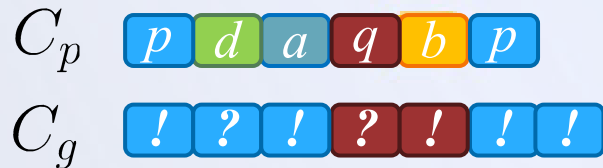
non-gaps

Encoding Event Sequences

Data D : $a \ b \ d \ c \ a \ d \ b \ a \ a \ b \ c$



Encoding 2: using patterns



The length of a gap code ? for pattern X

$$L(\text{?}) = -\log(p(\text{?} \mid p))$$

and analogue for non-gap codes !

Encoding Event Sequences

By which, the encoded size of D given CT and C is

$$L(D \mid CT) = L(C_p \mid CT) + L(C_g \mid CT)$$

...skipping the details of $L(CT \mid C)$...

Then, our goal is to minimise

$$L(CT, D) = L(CT \mid C) + L(D \mid CT)$$

Summarising Event Sequences

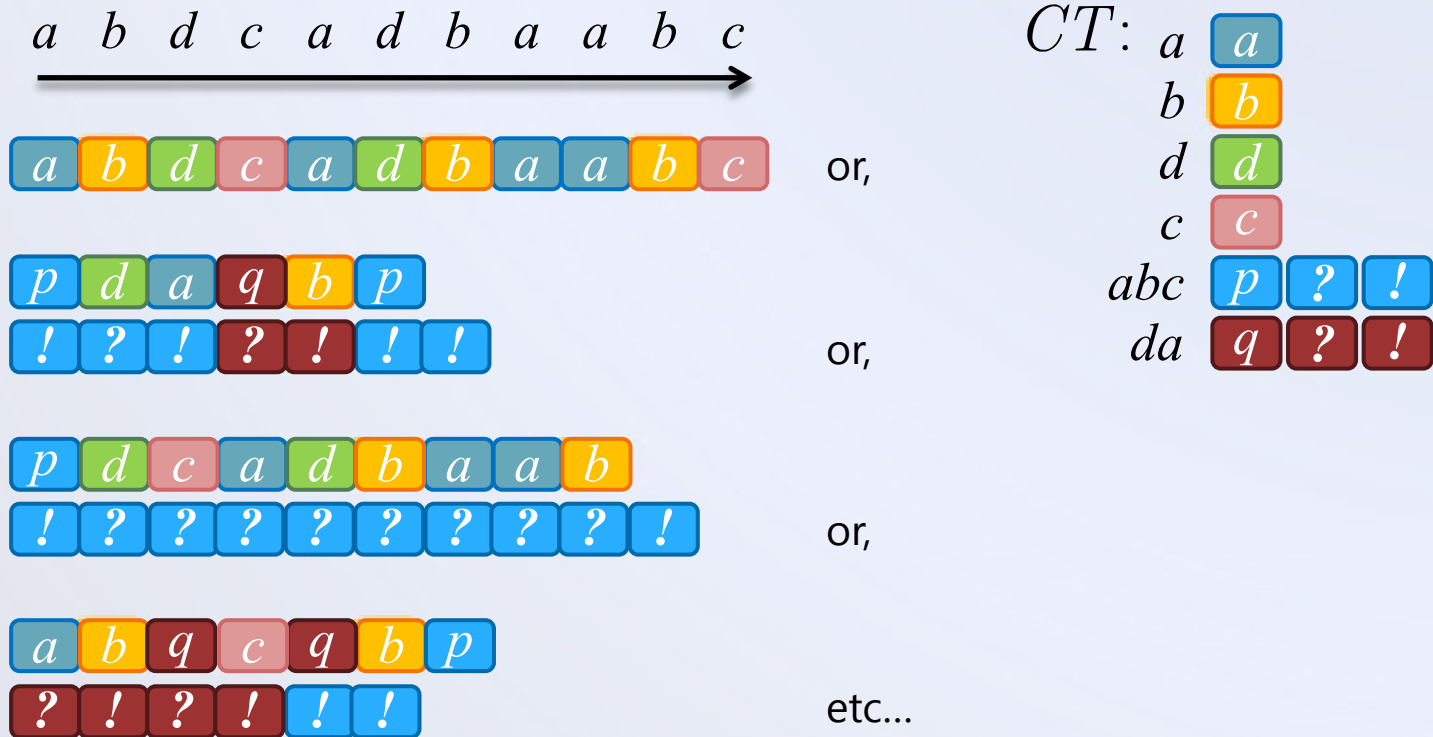
We want to find good summaries.

Three important questions

1. how do we **score** a summary?
2. how do we **describe** a sequence given a pattern set?
3. how do we **find** good pattern sets?

How to Cover your String

There are many ways C to describe a sequence given a set of patterns. We are after the **optimum**.



How to Cover your String

There are many ways C to describe a sequence given a set of patterns. We are after the **optimum**.

1. if we fix the **cover**, we can obtain the optimal code lengths
2. if we fix the **code lengths**, we can obtain the optimal cover by dynamic programming

We alternate these steps until **convergence**

How to Cover your String

压缩

There

given a

1. i

2. i

ng

We

How to Cover your String

sqs

Summarising event seQuenceS

Summarising Event Sequences

We want to find good summaries.

Three important questions

1. how do we **score** a summary?
2. how do we **describe** a sequence given a pattern set?
3. how do we **find** good pattern sets?

Mining Code Tables

There are very many possible pattern sets.

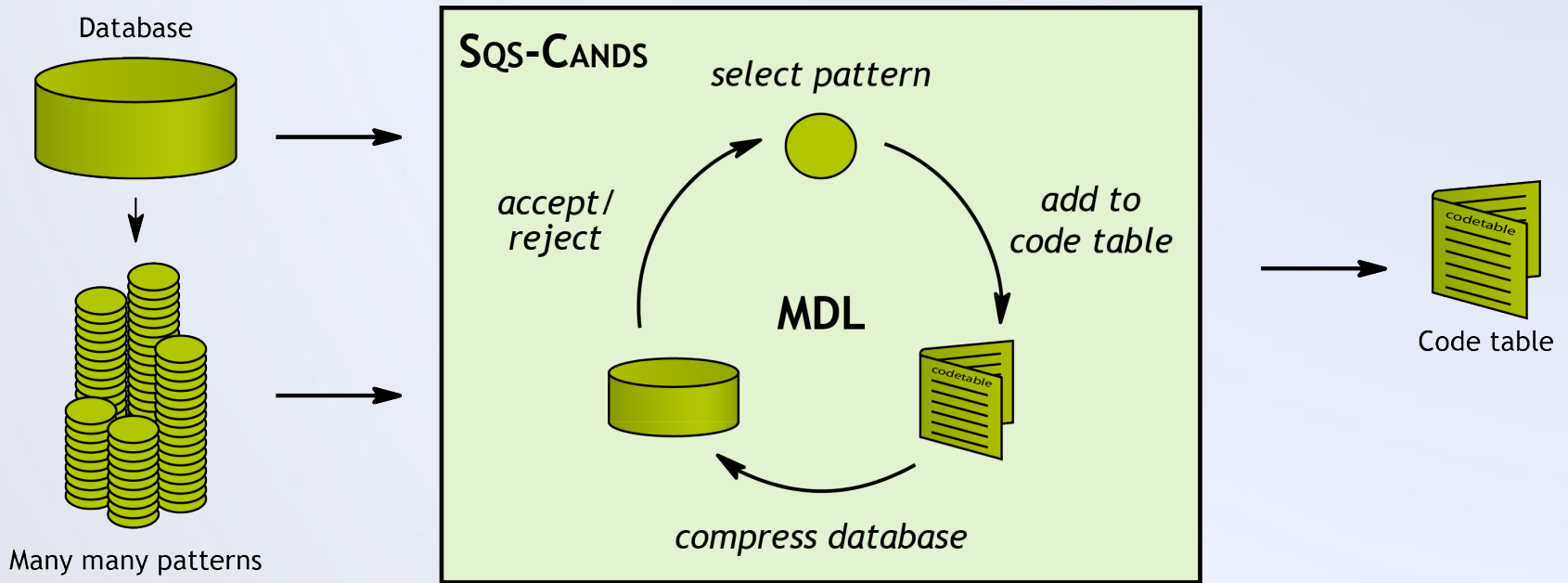
We are after the **optimum**

However, the search space is huge, complex, and does **not** exhibit trivial structure

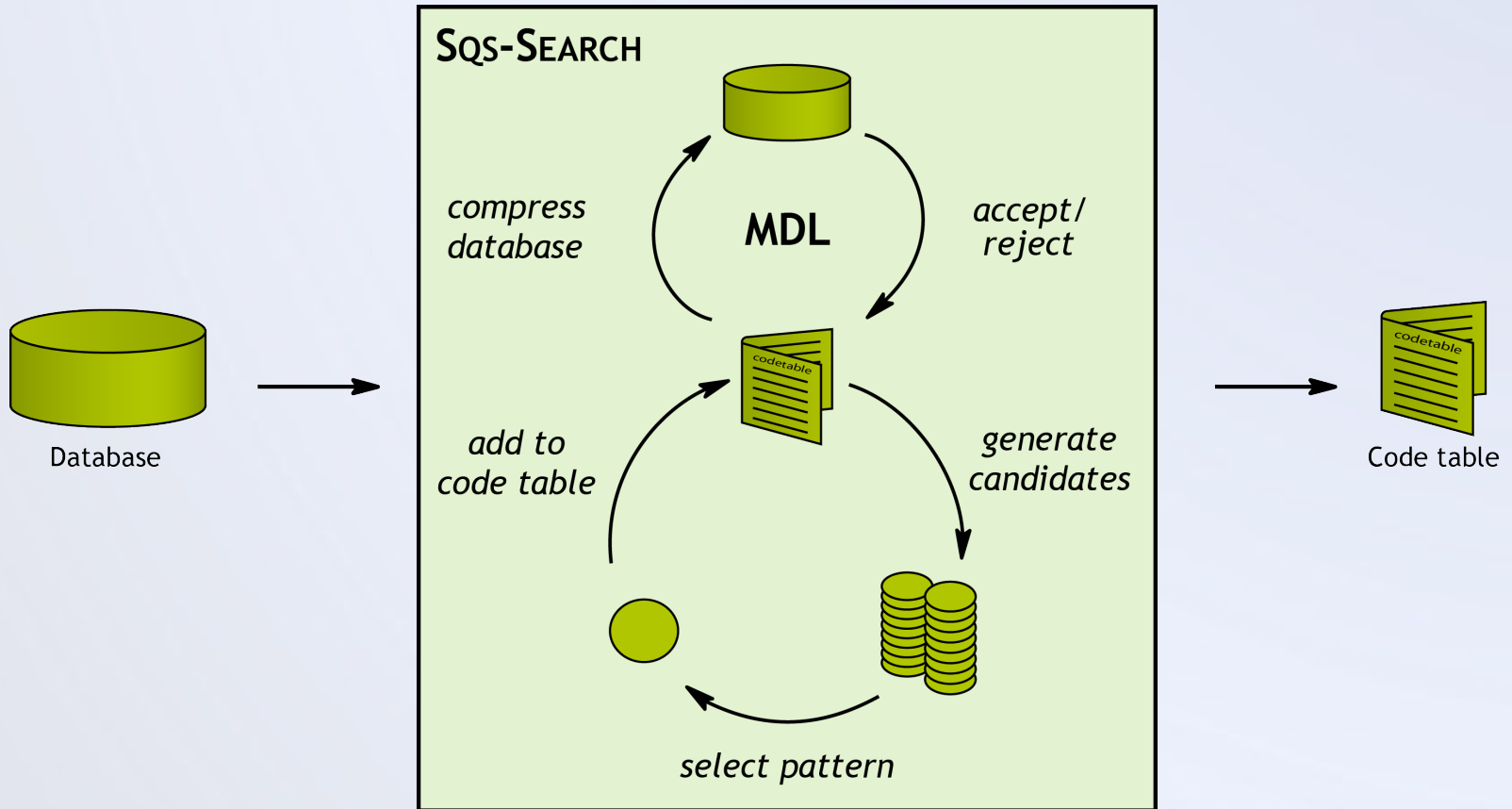
We propose two algorithms for mining code tables

- **SQS-CANDS** filters ordered lists of pre-mined candidates
- **SQS-SEARCH** mines good code tables directly from data

SQS-CANDIDATES



SQS-SEARCH



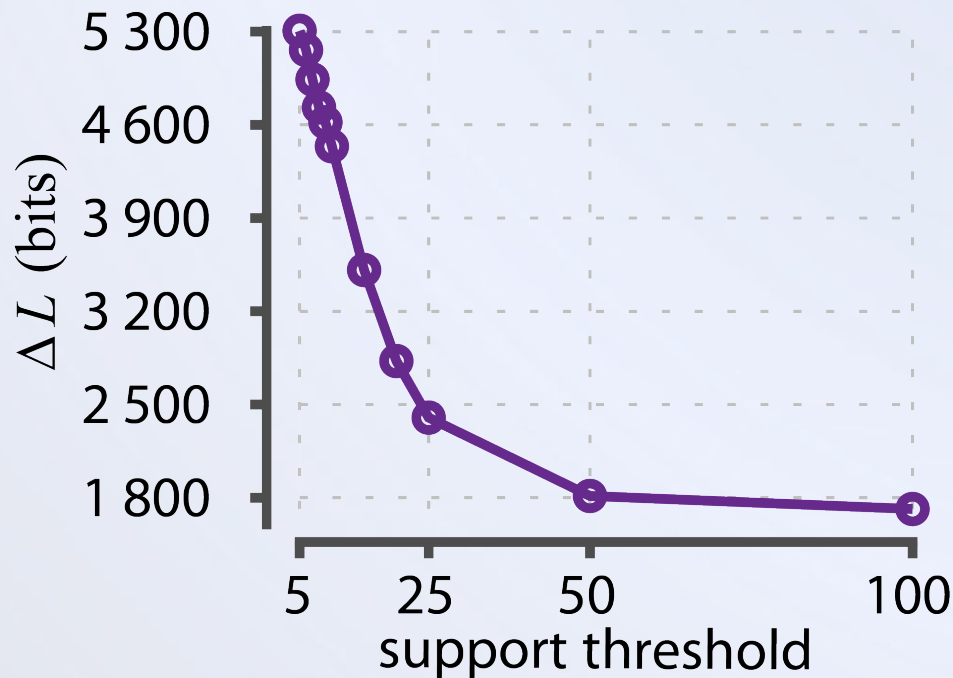
Experiments

- **synthetic data** random HMM ✓ no structure found
- **real data** text data ✓ structure recovered for interpretation

| | $ \Omega $ | $ D $ | SQS-CANDS | | SQS-SEARCH | ΔL |
|-----------|------------|-------|-----------|-------|------------|------------|
| | | | $ F $ | $ P $ | $ P $ | |
| Addresses | 5 295 | 56 | 15 506 | 138 | 155 | 5k |
| JMLR | 3 846 | 788 | 40 879 | 563 | 580 | 30k |
| Moby Dick | 10 277 | 1 | 22 559 | 215 | 231 | 10k |

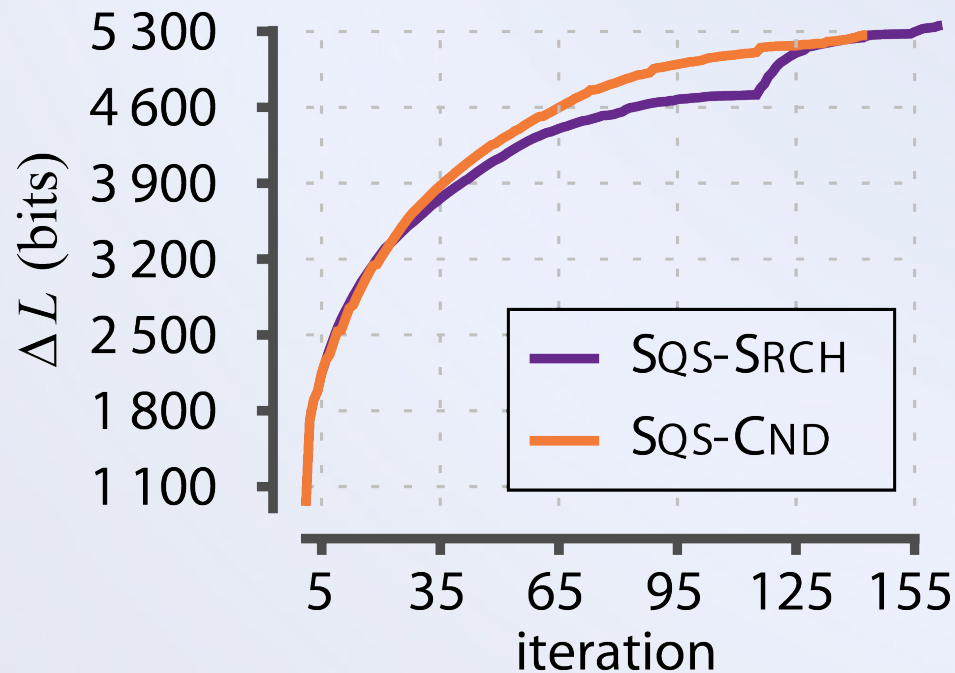
SQS-CANDIDATES

Compression improves with richer candidate sets
i.e. lower support thresholds



Optimising our Score

Both strategies show good convergence
SQS-SEARCH dips due to batch-wise search



Selected Results

JMLR

support vector machine
machine learning
state [of the] art
data set
Bayesian network

PRES. ADDRESSES

unit[ed] state[s]
public econ. expenditur
take oath
equal right
exercis power

Conclusions

Mining informative sets of patterns

- is an important aspect of exploratory data mining

SQS approximates the ideal for serial episodes

- SQS-CANDS filters a pre-mined candidate list
- SQS-SEARCH mines good code tables directly from data

Future work includes

- richer data and pattern types
- applying SQS in real-world settings

Thank you!

Mining informative sets of patterns

- is an important aspect of exploratory data mining

SQS approximates the ideal for serial episodes

- SQS-CANDS filters a pre-mined candidate list
- SQS-SEARCH mines good code tables directly from data

Future work includes

- richer data and pattern types
- applying SQS in real-world settings