

Symbolic Representations

EE382V Activity Sensing and Recognition

Today

Discuss findings from audio sensing activity

Symbolic representations for time series

A new way of thinking about activity recognition problems

Symbolic Aggregate Approximation (SAX)

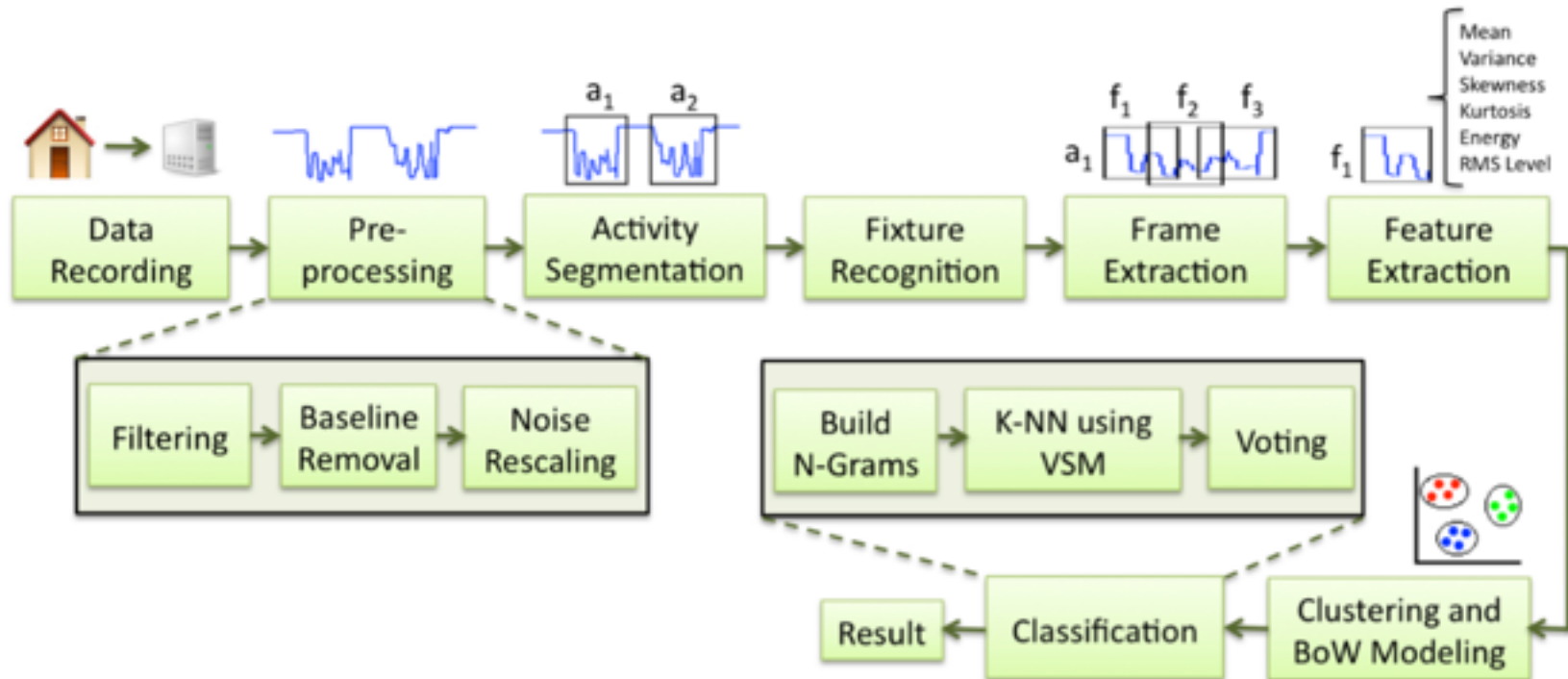
Random Projections

Audio Sensing Lab

Discuss findings from audio sensing activity (by group)

- How did you extract audio features?
- Which features did you choose?
- Which ML algorithms did you try training your classifier with?
- Which resulted in best performance?
- Which technique did you use to evaluate your classifier?
- How did your classifier perform with the test files?
- What did you find most challenging about this exercise?

Vector Space Model for Activity Recognition



Vector Space Model

Documents and queries are represented as vectors.

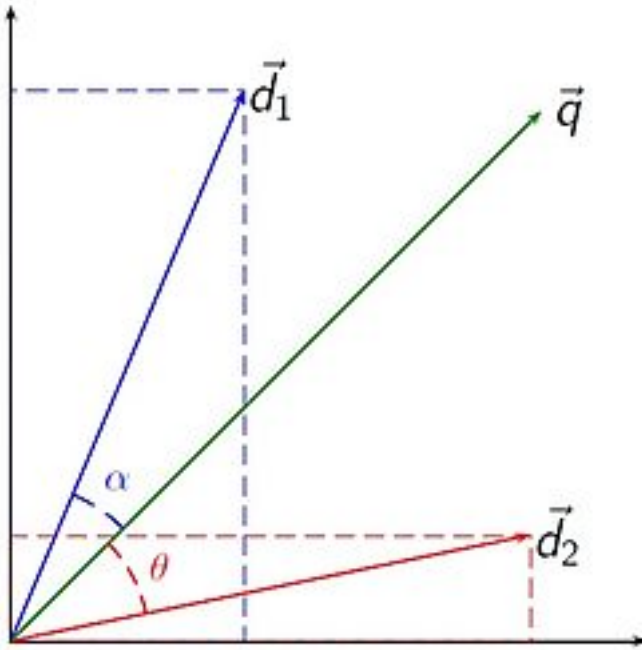
$$\begin{aligned} d_j &= (w_{1,j}, w_{2,j}, \dots, w_{t,j}) \\ q &= (w_{1,q}, w_{2,q}, \dots, w_{n,q}) \end{aligned} \quad \text{(term) weights,}$$

If a term occurs in the document, its value in the vector is non-zero

$$w_{t,d} = \text{tf}_{t,d} \cdot \log \frac{|D|}{|\{d' \in D \mid t \in d'\}|}$$

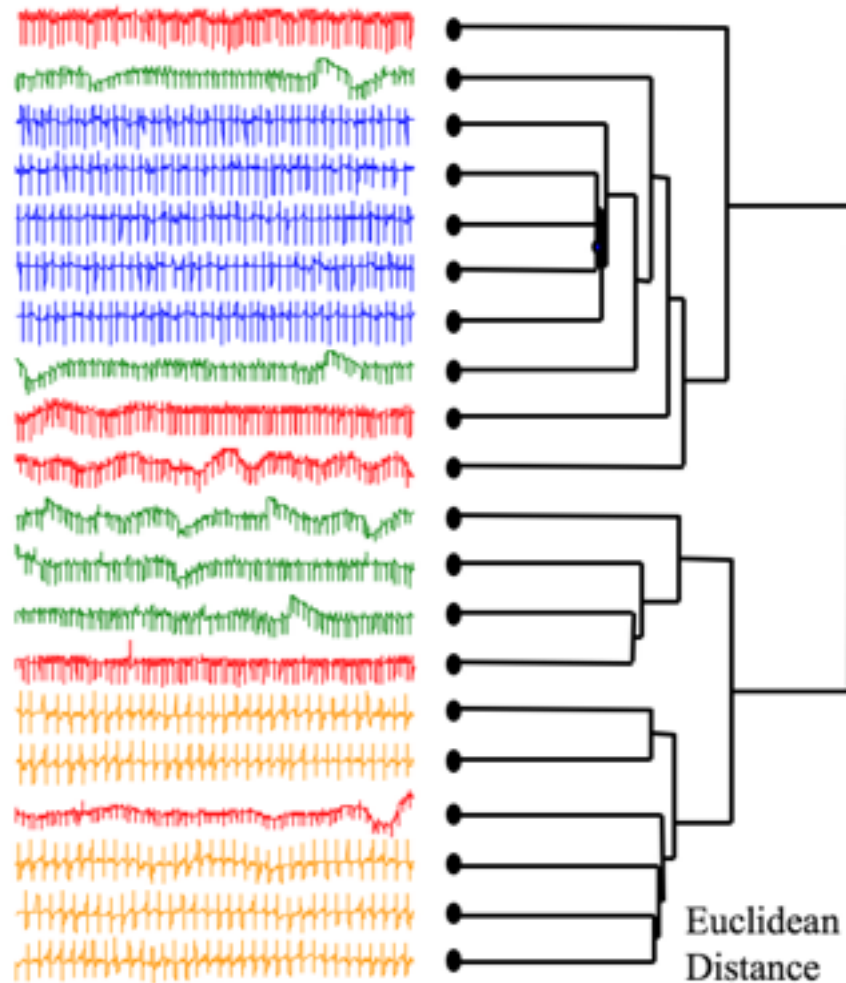
tf-idf weights

Vector Space Model



Relevance rankings of documents in a keyword search can be calculated by comparing the deviation of angles between each document vector and the original query vector

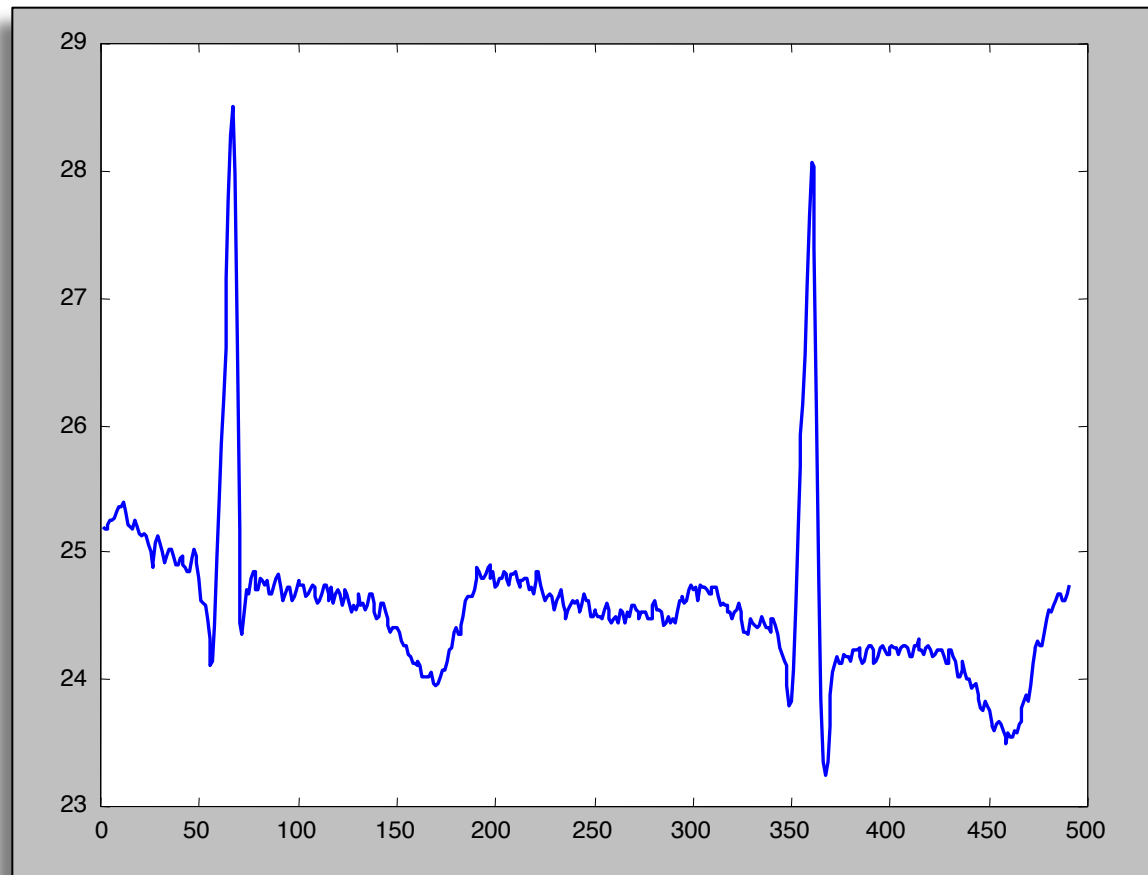
$$\cos \theta = \frac{\mathbf{d}_2 \cdot \mathbf{q}}{\|\mathbf{d}_2\| \|\mathbf{q}\|}$$



Slides thanks to Eamon Keogh and Jessica Lin
University of California Riverside

What are Time Series?

A time series is a collection of observations made sequentially in time.



25.1750
25.2250
25.2500
25.2500
25.2750
25.3250
25.3500
25.3500
25.4000
25.4000
25.3250
25.2250
25.2000
25.1750

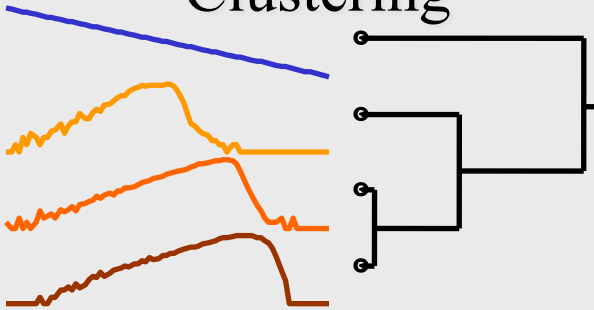
••

••

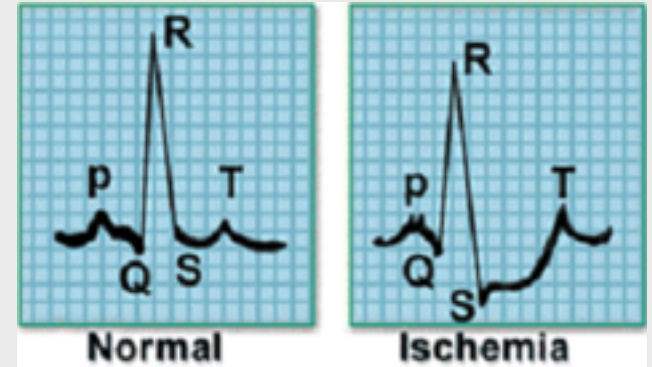
24.6250
24.6750
24.6750
24.6250
24.6250
24.6250
24.6750
24.7500

What do we want to do with the time series data?

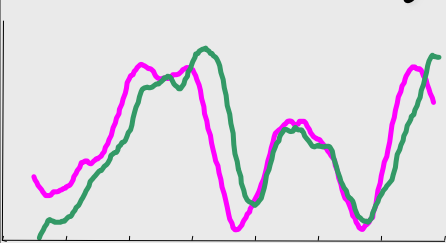
Clustering



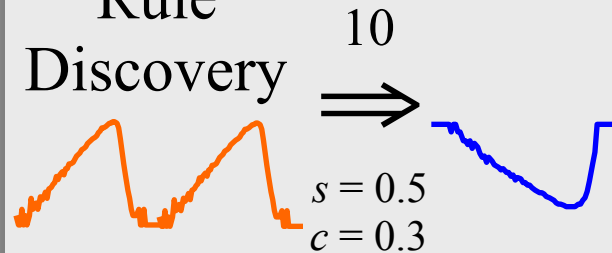
Classification



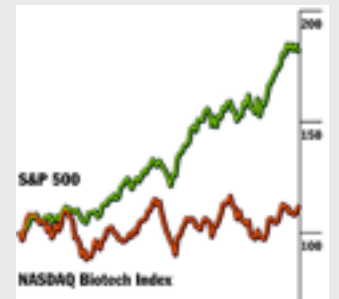
Motif Discovery



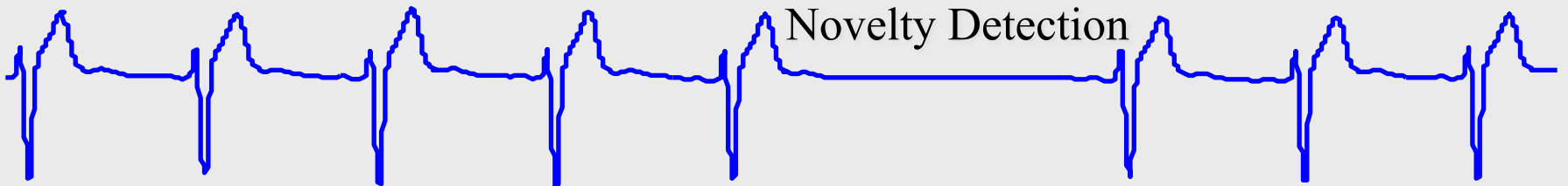
Rule
Discovery



Query by
Content

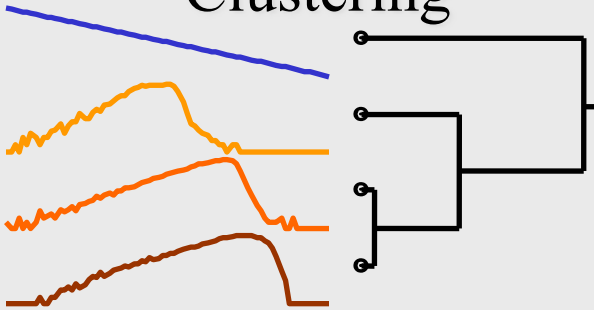


Novelty Detection

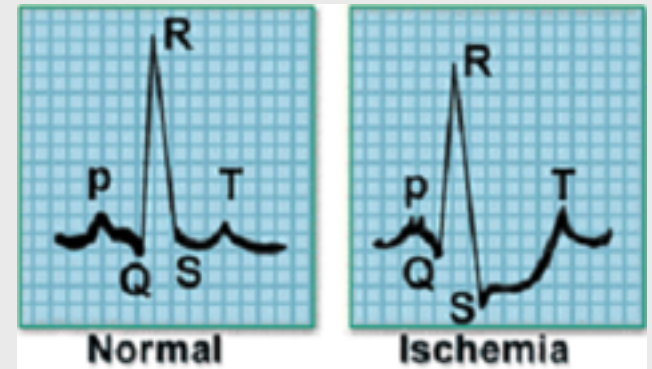


All these problems require **similarity** matching

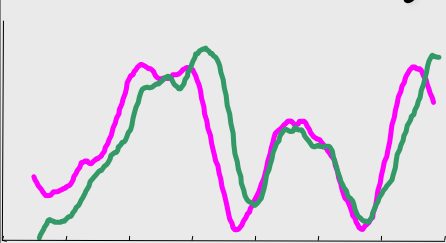
Clustering



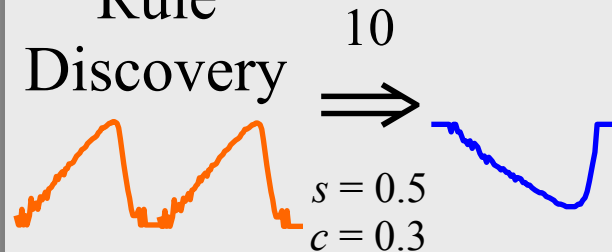
Classification



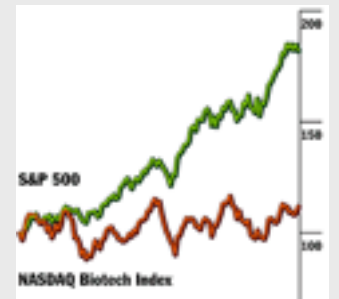
Motif Discovery



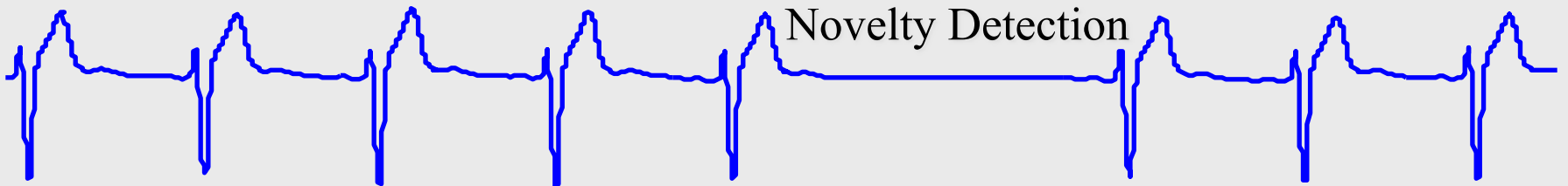
Rule
Discovery



Query by
Content



Novelty Detection



Euclidean Distance Metric

Given two time series

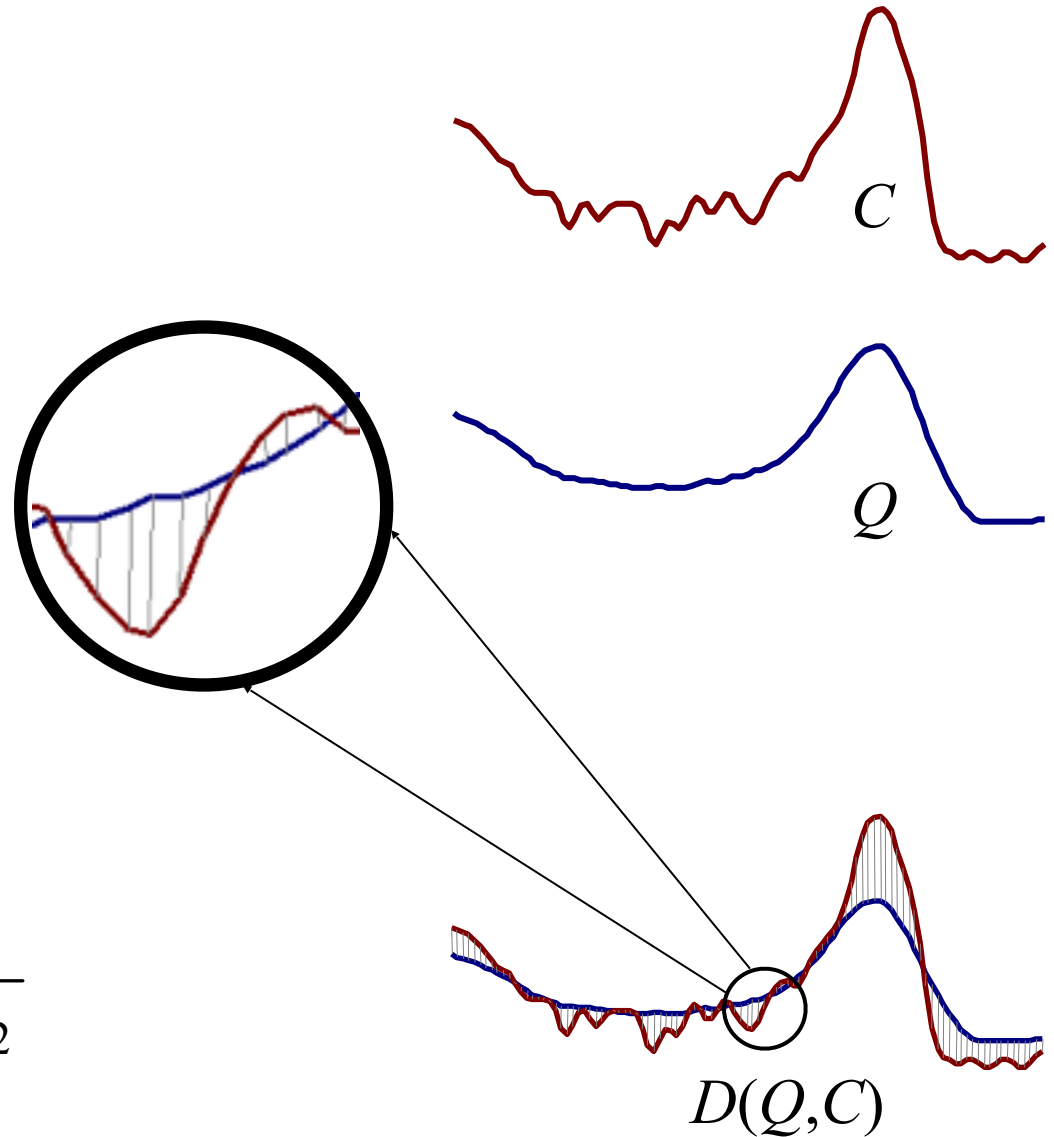
$$Q = q_1 \dots q_n$$

and

$$C = c_1 \dots c_n$$

their Euclidean distance is defined as:

$$D(Q, C) \equiv \sqrt{\sum_{i=1}^n (q_i - c_i)^2}$$



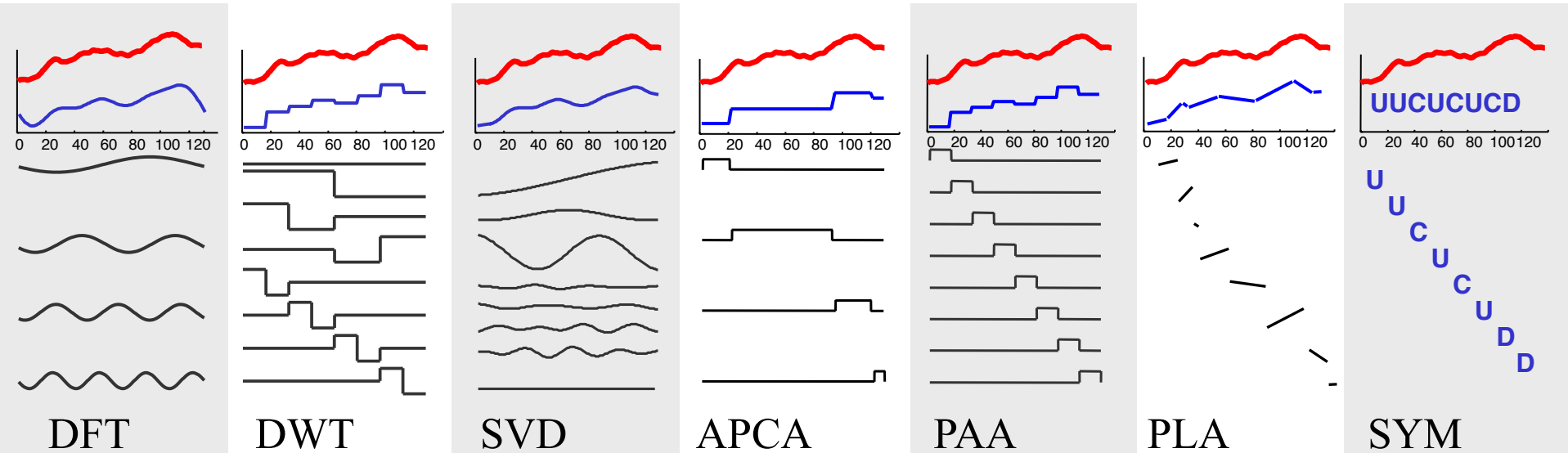
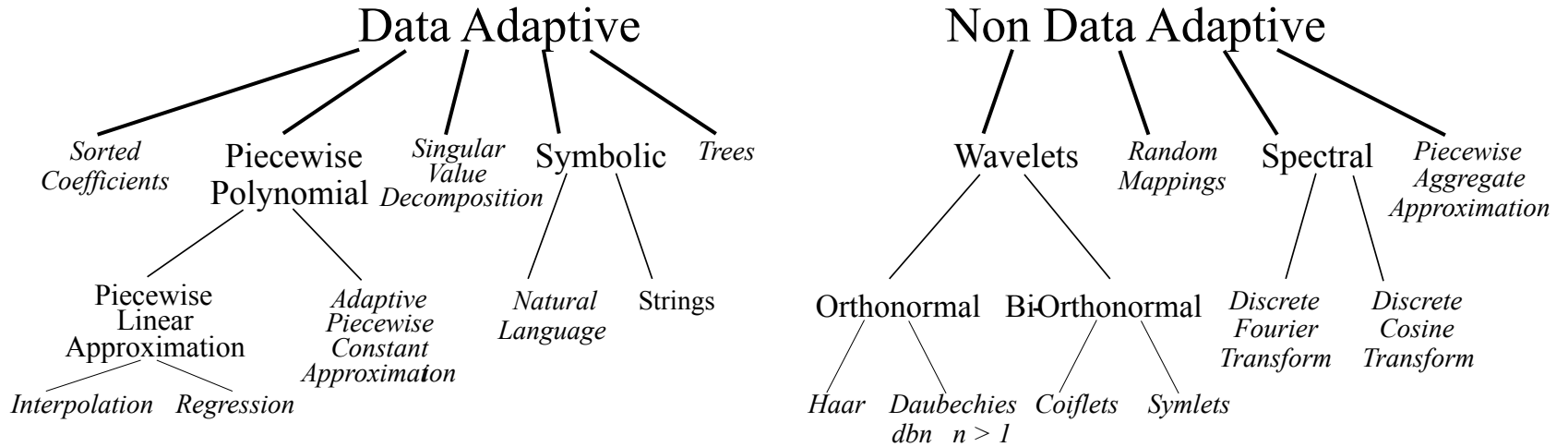
The Generic Data Mining Algorithm

- Create an *approximation* of the data, which will fit in main memory, yet retains the essential features of interest
- Approximately solve the problem at hand in main memory
- Make (hopefully very few) accesses to the original data on disk to confirm the solution obtained in Step 2, or to modify the solution so it agrees with the solution we would have obtained on the original data

But which approximation
should we use?



Time Series Representations



The Generic Data Mining Algorithm (revisited)

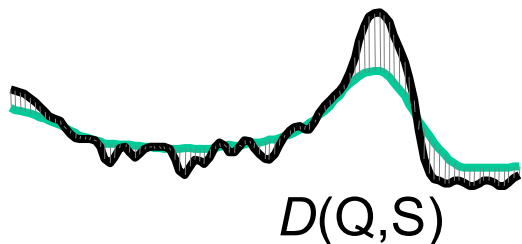
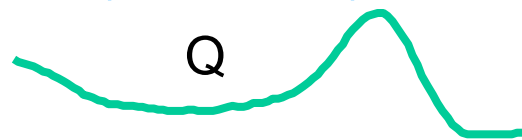
- Create an *approximation* of the data, which will fit in main memory, yet retains the essential features of interest
- Approximately solve the problem at hand in main memory
- Make (hopefully very few) accesses to the original data on disk to confirm the solution obtained in Step 2, or to modify the solution so it agrees with the solution we would have obtained on the original data

This only works if the
approximation allows
lower bounding



What is lower bounding?

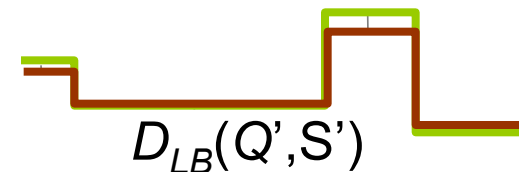
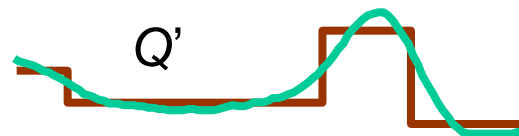
Exact (Euclidean) distance $D(Q,S)$



$D(Q,S)$

$$\equiv \sqrt{\sum_{i=1}^n (q_i - s_i)^2}$$

Lower bounding distance $D_{LB}(Q,S)$



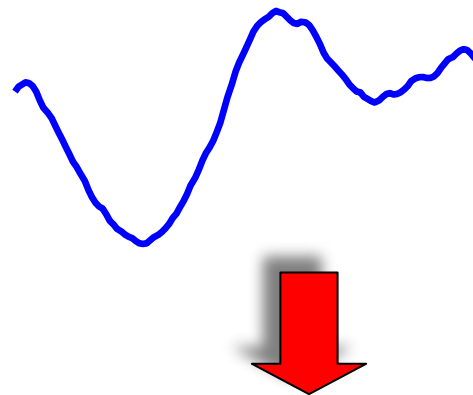
$D_{LB}(Q',S')$

$$\equiv \sqrt{\sum_{i=1}^M (sr_i - sr_{i-1})(qv_i - sv_i)^2}$$

Lower bounding means that for all Q and S, we have...

$$D_{LB}(Q',S') \leq D(Q,S)$$

Symbolic Aggregate ApproXimation (SAX)

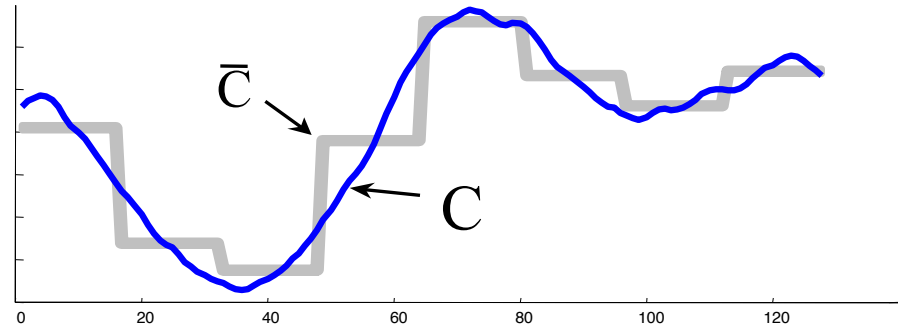


baabccbc

- Lower bounding of Euclidean distance
- Dimensionality Reduction
- Numerosity Reduction

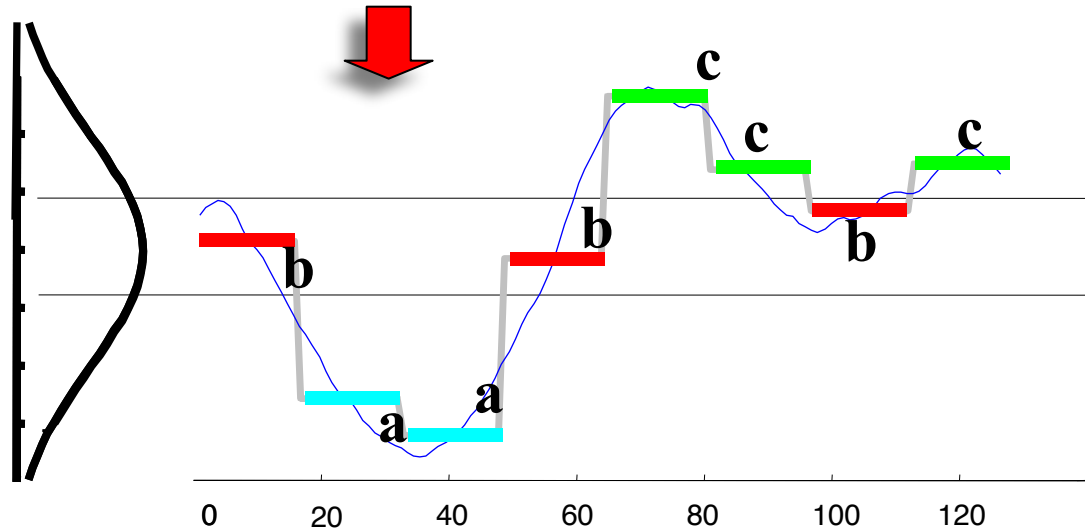


How do we obtain SAX?

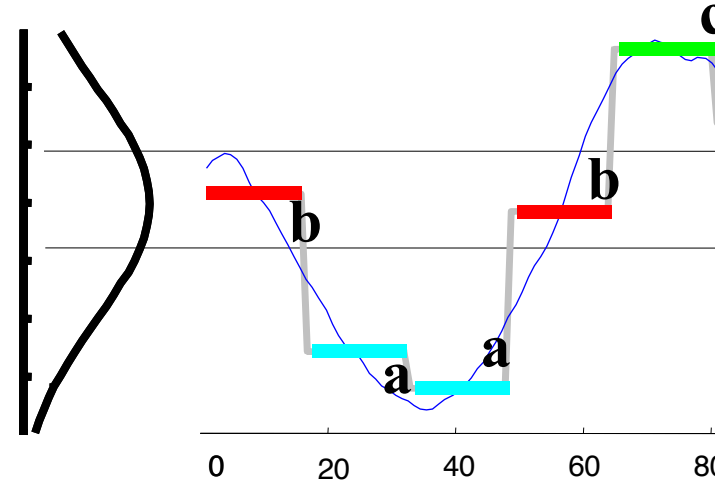


First convert the time series to PAA representation, then convert the PAA to symbols

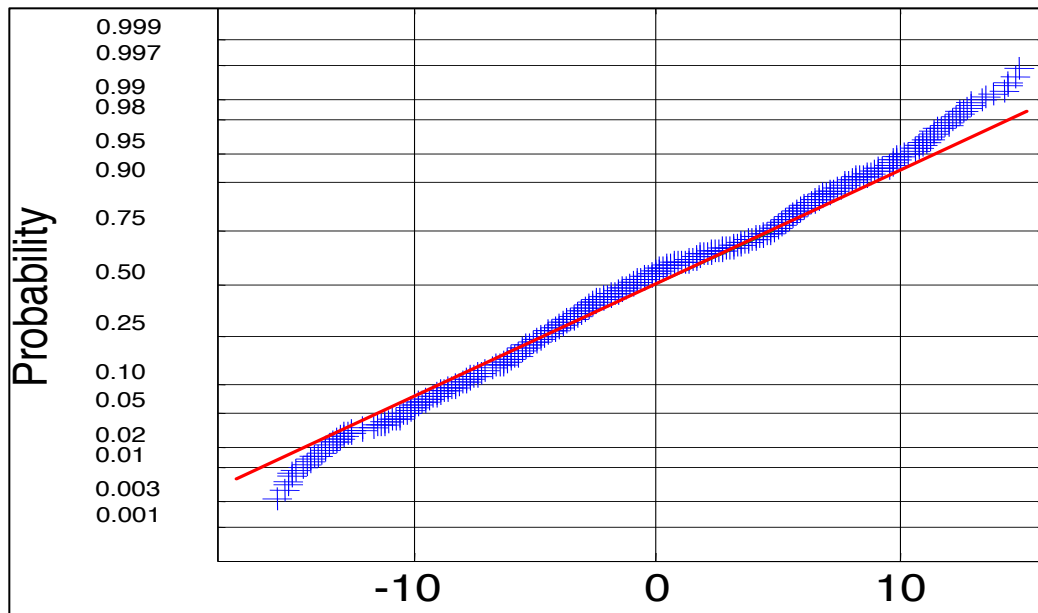
It take linear time



baabccbc



Time series subsequences tend to have a highly Gaussian distribution

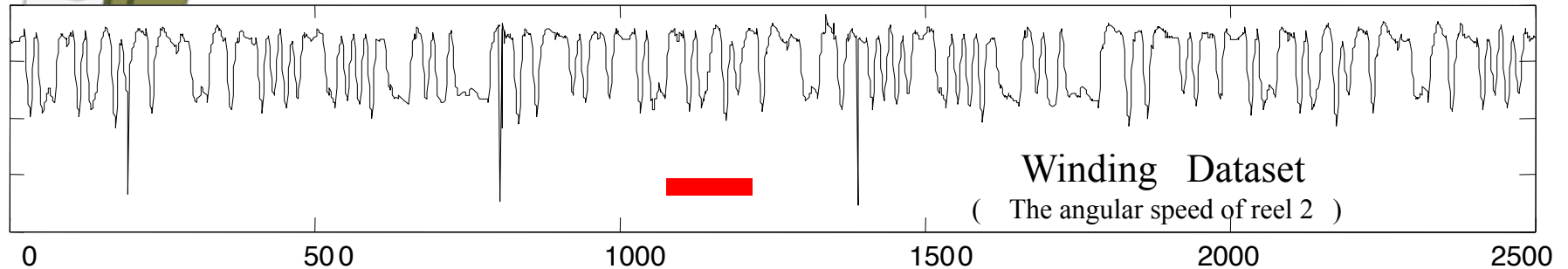


A normal probability plot of the (cumulative) distribution of values from subsequences of length 128.

Why a Gaussian?



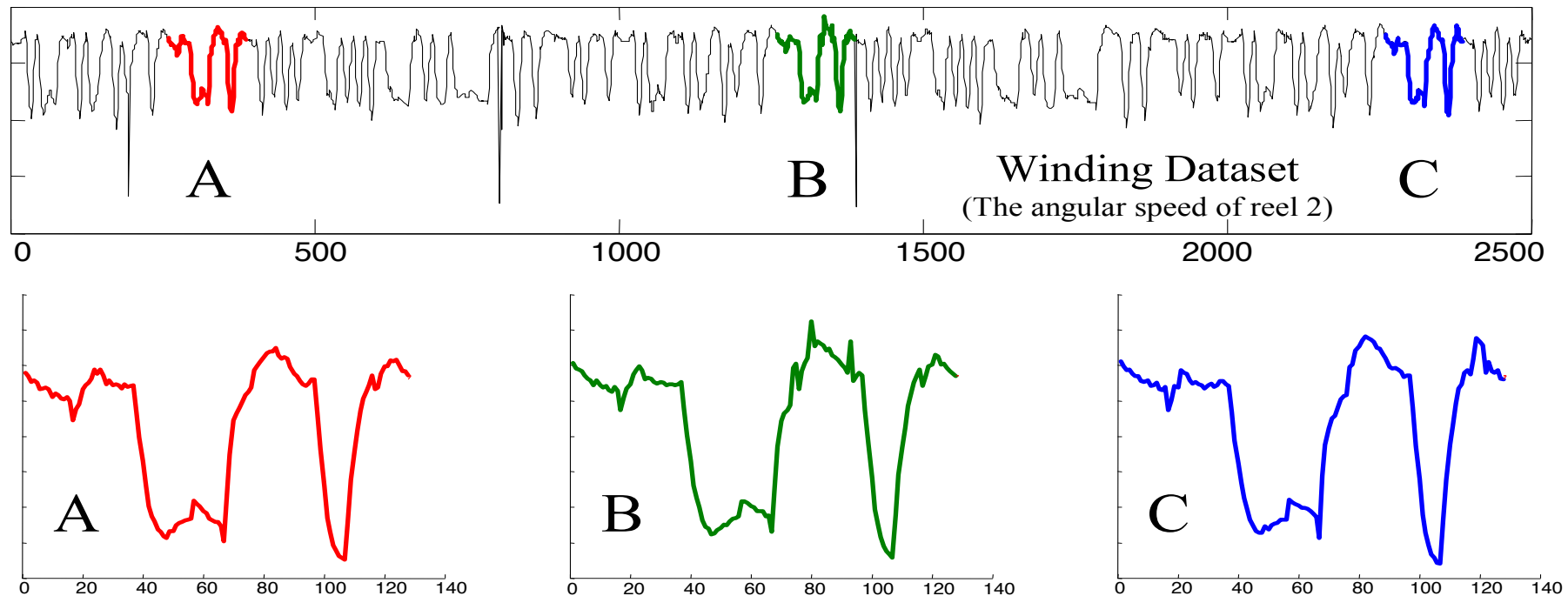
SAX allows Motif Discovery!



Informally, motifs are reoccurring patterns...

Motif Discovery

To find these 3 motifs would require about 6,250,000 calls to the Euclidean distance function.



OK, we can define motifs, but how do we find them?

The obvious brute force search algorithm is just too slow...

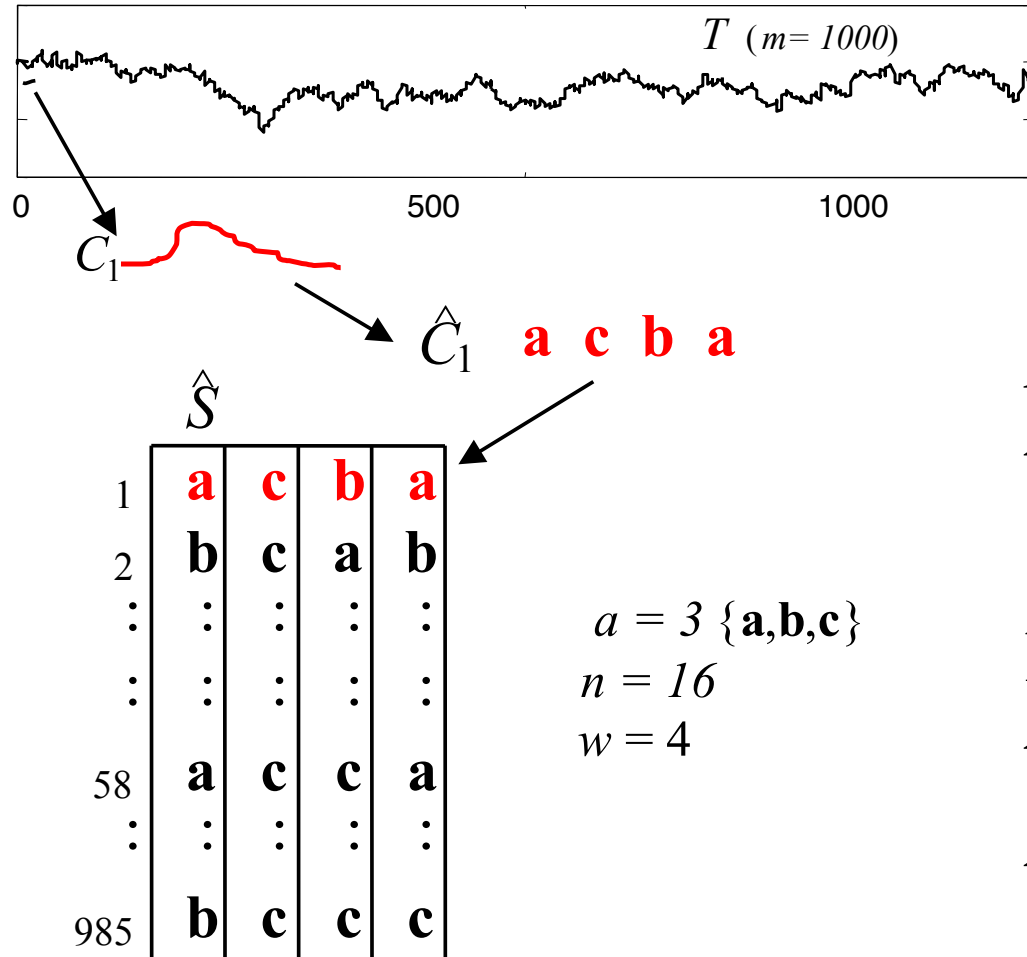
Our algorithm is based on a *hot* idea from bioinformatics, *random projection** and the fact that SAX allows use to lower bound discrete representations of time series.

* J Buhler and M Tompa. *Finding motifs using random projections*. In RECOMB'01. 2001.



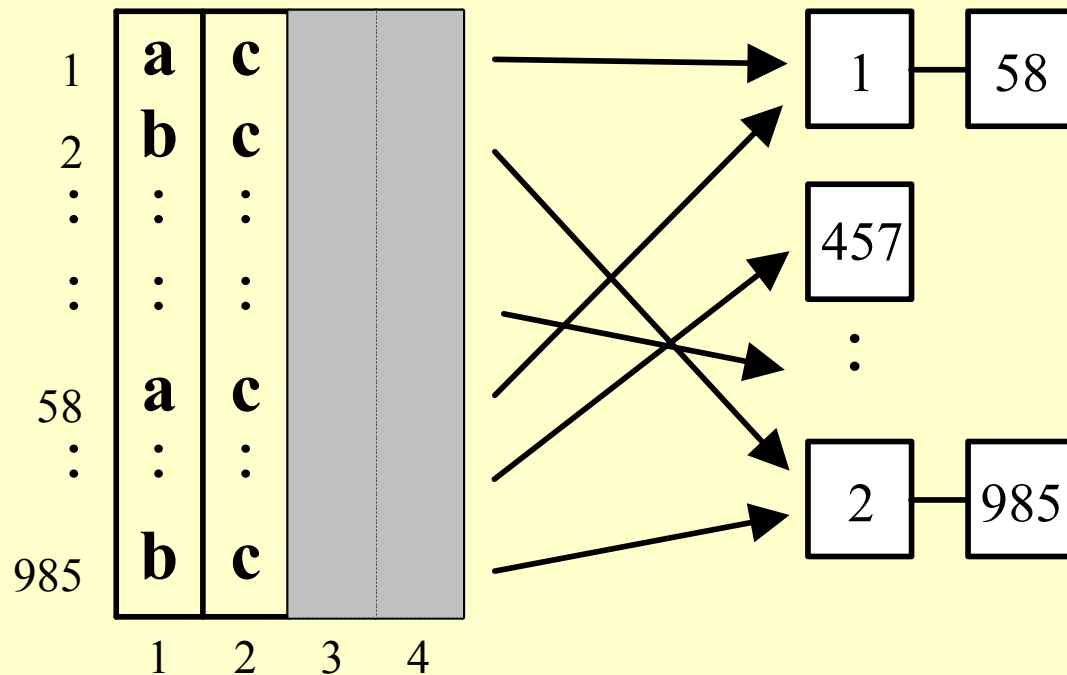
A simple worked example of our motif discovery algorithm

The next 4 slides



Assume that we have a time series T of length 1,000, and a motif of length 16, which occurs twice, at time T_1 and time T_{58} .

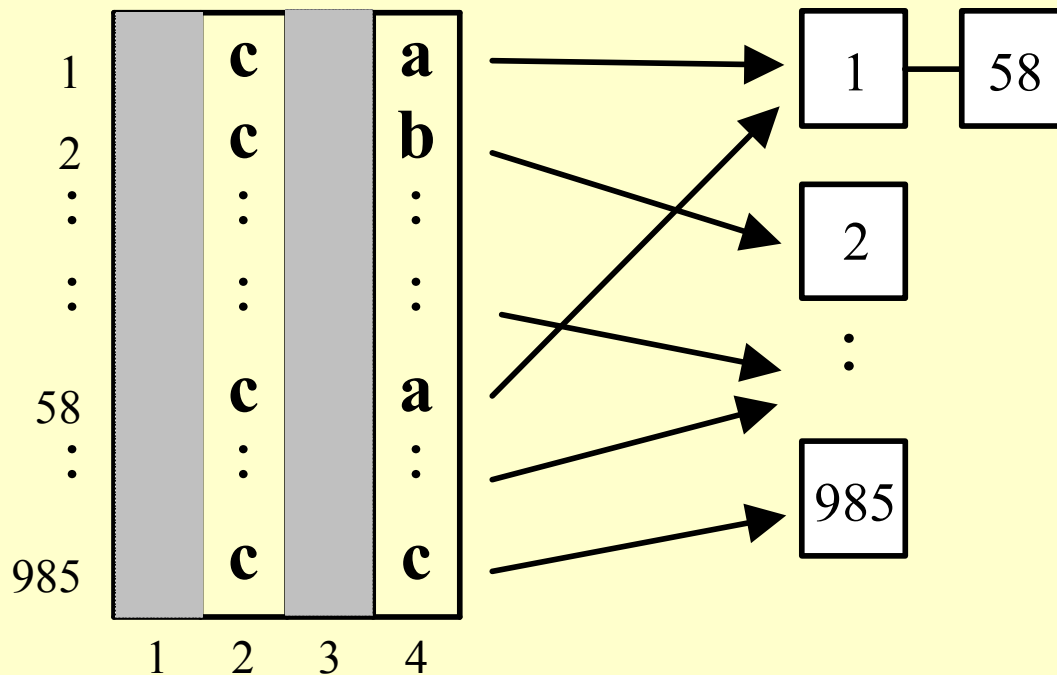
A mask $\{1,2\}$ was randomly chosen, so the values in columns $\{1,2\}$ were used to project matrix into buckets.



Collisions are recorded by incrementing the appropriate location in the collision matrix

1					
2					
:					
58	1				
:					
985		1			
	1				
	2	:	58	:	985

A mask $\{2,4\}$ was randomly chosen, so the values in columns $\{2,4\}$ were used to project matrix into buckets.



Once again, collisions are recorded by incrementing the appropriate location in the collision matrix

1						
2						
:						
58	2					
:						
985		1				
	1	2	:	58	:	985

1	a	c	b	a
2	b	c	a	b
:	:	:	:	:
:	:	:	:	:
58	a	c	c	a
:	:	:	:	:
985	b	c	c	c

1					
2	2				
:	1	3			
58	27	2	1		
:	3	2	2	1	
985	0	1	2	1	3
	1	2	:	58	:
					985

Coming up...

Project Progress Report Due Today at 11:59PM

Panel of Experts (Thursday)

AR Models from Data and Human Computation

Panel of Experts (Next Tuesday)

Environmental Sensing: Objects and Simple Sensors