# Mining Time Series Data 1

# Acknowledgements

Some slides are provided by

- Eamonn Keogh (various time series data mining tutorials)
- Michail Vlachos (ICDM 2006 tutorial)

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

# •• 24.6250 24.6750 24.6250 24.6250 24.6250 24.6250 24.6750 24.7500

...

# What are Time Series?

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





# Time Series are Ubiquitous! I

People measure things...

- Their blood pressure
- George Bush's popularity rating
- The annual rainfall in Seattle
- The value of their Google stock

... and things change over time...

Thus time series occur in virtually every medical, scientific and businesses domain

# Time Series are Ubiquitous! II

A random sample of 4,000 graphics from 15 of the world's newspapers published from 1974 to 1989 found that more than 75% of all graphics were time series (Tufte, 1992).



# Shapes

#### Recognize type of leaf based on its shape











Ulmus carpinifolia

Acer platanoides

Salix fragilis

Tilia

Quercus robur

### Convert perimeter into a sequence of values



ICDM 2006 Tutorial Michail Vlachos



# Image Histograms



# Applications (Image Matching)

**Cluster 1** 

Many types of data can be converted to time-series

Image



Color Histogram





# Applications (Motion Capture)

### Motion-Capture (MOCAP) Data (Movies, Games)

MOCAP data...

...my precious...

- Track position of several joints over time
- 3\*17 joints = **51 parameters per frame**







# Applications (Video)

Video-tracking / Surveillance

- Visual tracking of body features (2D time-series)
  - Sign Language recognition (3D time-series)





Text data, may best be thought of as time series...

# The local frequency of words in the Bible



# **Text Data As Time Series II**



Regional interest 2



India	100
South Korea	66
Singapore	63
United States	40
Hong Kong	40
Taiwan	38
South Africa	31

Region City

# **Text Data As Time Series III**



### Handwriting data, may best be thought of as time series...

Setters in 1158. it and to prevent this advantageous Commerce from fuffing in its infancy by the simister views of defigning, subject men, of the different Provinces - I have the conceive it absolutely necessary, that Commissioners from each of the almies be appointed, to regulate the move of that Irade, and fix it on such a bafis that, all the attempts of one Colony anming and shereby weakening and diminishing the gentral fyster, might he pustines " Do effect which the general would ( I anoy ) cheapily give his aid Micho none can entertain a higher Sense of the greate importance of main taining a Post upon the this then mys elf, get under the unhappy ai curretances that my Regiment is Iwould by no means have agreed to have any part of it there, hav not the Gos queer an express aver for it Sent voured to show that the Kings Dorops ought to galenifor togetes while total these hav in surpland my fable Stuations sela tasti Da on is profile any the theo with the language egts pull ! and, if the First U. Requirent

**George Washington Manuscript** 



George Washington 1732-1799



### Brain scans (3D voxels), may best be thought of as time series..







Works with 3D glasses!







Wang, Kontos, Li and Megalooikonomou ICASSP 2004

# Why is Working With Time Series so Difficult? Part I

### Answer: How do we work with very large databases?

- 1 Hour of ECG data: 1 GB
- Twitter: over 400 million tweets per day (in 2013, compared to 340M in 2012, and 200M in 2011)
- Astronomy Databases:
  - LSST (Large Synoptic Survey Telescope) project: 20 PB science data & 100 PB image archive
- Satellite data
- Sensor data
- Biological data

Since most of the data lives on disk (or tape), we need a representation of the data we can efficiently manipulate.

# Why is Working With Time Series so Difficult? Part II

Answer: We are dealing with subjectivity



The definition of similarity depends on the user, the domain and the task at hand. We need to be able to handle this subjectivity.

# Why is working with time series so difficult? Part III

Answer: Miscellaneous data handling problems.

- Differing data formats.
- Differing sampling rates.
- Noise, missing values, etc.

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



### All these problems require similarity matching



### A simple motivation



You go to the doctor because of chest pains. Your ECG looks strange...

You doctor wants to search a database to find **similar** ECGs, in the hope that they will offer clues about your condition...

Two questions:

- How do we define similar?
- How do we search quickly?

# **Basic Time-Series Matching Problem**



### Linear Scan:

query

**Objective:** Compare the query with all sequences in DB and return the k most similar sequences to the query.

Database with time-series:

- Medical sequences
- Images, etc

Sequence Length:100-1000pts DB Size: 1 TByte



# What other problems can we solve?

Clustering: "Place time-series into <u>'similar'</u> groups"



Classification: "To which group is a time-series most <u>'similar'</u> to?"



# Hierarchical Clustering

- Very generic & powerful tool
- Provides visual data grouping

Pairwise distances



- 1. Merge objects with smallest distance
- 2. Re-evaluate distances
- 3. Repeat process



Z = linkage(D); H = dendrogram(Z);

# Partitional Clustering

- Faster than hierarchical clustering
- Typically provides suboptimal solutions (local minima)
- Not good performance for high dimensions

### K-Means Algorithm:

- 1. Initialize k clusters (k specified by user) randomly.
- 2. Repeat until convergence
  - 1. Assign each object to the nearest cluster center.
  - 2. Re-estimate cluster centers.



# K-Means Demo



# K-Means Clustering for Time-Series

- So how is kMeans applied for Time-Series that are highdimensional?
- Perform kMeans on a compressed dimensionality



# Classification

# Typically classification can be made easier if we have clustered the objects



So, query Q is more similar to class B

# Nearest Neighbor Classification

We need not perform clustering before classification. We can classify an object based on the class majority of its nearest neighbors/matches.





# What is Similarity? The quality or state of being similar; likeness; resemblance; as, a similarity of features. Webster's Dictionary



Similarity is hard to define, but... "We know it when we see it"

The real meaning of similarity is a philosophical question.

We will take a more pragmatic approach.

# Notion of Similarity I

 Solution to any time-series problem, boils down to a proper definition of <u>\*similarity\*</u>



Similarity is always <u>subjective</u>. (*i.e. it depends on the application*)

# Notion of Similarity II

#### Similarity depends on the *features* we consider

(i.e. how we will describe or compress the sequences)







rockbass.gif 790 x 428 pixels - 201k www.dec.state.ny.us/.../ fishspecs/rockbass.gif



bass5.jpg 600 x 800 pixels - 120k www.danielbuttner.com/ bass4sale/bass5.jpg



**bass**-Irgmouth.jpg 722 x 432 pixels - 32k agrino.org/.../freshfish/ bass-Irgmouth.jpg



m082\_anejo-**bass**.jpg 559 x 795 pixels - 71k www.personal.rdg.ac.uk/.../ m082\_anejo-bass.jpg












### Metric and Non-metric Distance Functions



#### **Properties**

**Positivity and Constancy:** 

 $d(x,y) \ge 0$  and d(x,y)=0, if x=y

<u>Symmetry:</u> d(x,y) = d(y,x)

<u>Triangle Inequality:</u>  $d(x,z) \le d(x,y) + d(y,z)$ 

*If <u>any</u> of these is not obeyed then the distance is a non-metric* 



## Triangle Inequality (Importance)

<u>Triangle Inequality:</u> d(x,z) ≤ d(x,y) + d(y,z)



 Assume:
 d(Q,bestMatch) = 20

 and
 d(Q,B) =150

Then, since d(A,B)=20

 $d(Q,A) \geq d(Q,B) - d(B,A)$ 

 $d(Q,A) \ge 150 - 20 = 130$ 

#### So we don't have to retrieve A from disk

	Α	В	С
A	0	20	110
В	20	0	90
С	110	90	0

### Non-Metric Distance Functions





- Matching flexibility
- Robustness to outliers
- Stretching in time/space
- Support for different sizes/lengths



 Speeding-up search can be difficult

# **Euclidean Distance Metric** Given two time series: $Q = q_1 \dots q_n$ $C = c_1 \dots c_n$ $D(Q,C) \equiv \sqrt{\sum_{i=1}^{n} (q_i - c_i)^2}$ About 80% of published work in data mining uses Euclidean distance D(Q,C)L2 = sqrt(sum((a-b).^2)); % return Euclidean distance

# **Optimizing the Euclidean Distance Calculation**

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

Euclidean distance and Squared Euclidean distance are equivalent in the sense that they return the same rankings, clusterings and classifications



### Euclidean Distance (Vectorization) A Matlab trick

**Question:** If I want to compare many sequences to each other do I have to use a for-loop?

Answer: No, one can use the following equation to perform matrix computations only...

||A-B|| = sqrt ( ||A||<sup>2</sup> + ||B||<sup>2</sup> - 2\*A.B )



aa=sum(a.\*a); bb=sum(b.\*b); ab=a'\*b; d = sqrt(repmat(aa',[1 size(bb,2)]) + repmat(bb,[size(aa,2) 1]) - 2\*ab);

#### Preprocessing the data before distance calculations

If we naively try to measure the distance between two "raw" time series, we may get very unintuitive results

This is because Euclidean distance is very sensitive to some "distortions" in the data. For most problems these distortions are not meaningful, and thus we can and should remove them

In the next few slides we will discuss the 4 most common distortions, and how to remove them

- Offset Translation
- Amplitude Scaling
- Linear Trend
- Noise

### **Transformation I: Offset Translation**





### **Transformation II: Amplitude Scaling**



 $Q = (Q - \operatorname{mean}(Q)) / \operatorname{std}(Q)$  $C = (C - \operatorname{mean}(C)) / \operatorname{std}(C)$ D(Q,C)

### **Transformation III: Linear Trend**



The intuition behind removing linear trend is...

Fit the best fitting straight line to the time series, then subtract that line from the time series.



Removed linear trend

Removed offset translation

Removed amplitude scaling

### **Transformation IIII: Noise**



The intuition behind removing noise is...

Average each datapoints value with its neighbors.

 $Q = \operatorname{smooth}(Q)$  $C = \operatorname{smooth}(C)$ D(Q,C)

# A Quick Experiment to Demonstrate the Utility of Preprocessing the Data



# **Summary of Preprocessing**

The "raw" time series may have distortions which we should remove before clustering, classification etc

> Of course, sometimes the distortions are the most interesting thing about the data, the above is only a general rule



We should keep in mind these problems as we consider the high level representations of time series which we will encounter later (DFT, Wavelets etc). Since these representations often allow us to handle distortions in elegant ways



# Euclidean distance cannot compensate for small distortions in time axis.



According to Euclidean distance B is more similar to A than to C

Solution: Allow for compression & decompression in time



Note: We will first see the utility of DTW, then see how it is calculated.



## Dynamic Time-Warping

First used in speech recognition for recognizing words spoken at different speeds

Same idea can work equally well for generic time-series data



### Dynamic Time-Warping (how does it work?)

The intuition is that we copy an element multiple times so as to achieve a better matching

Euclidean distance				
T1 = [1, 1, 2, 2]				
T2 = [1, 2, 2, 2]				

Warping distance T1 = [1, 1, 2, 2] 1 d = 0T2 = [1, 2, 2, 2]





**One-to-many non-linear alignment** 

## How is DTW Calculated? I

We create a matrix the size of |Q| by |C|, then fill it in with the distance between every pair of point in our two time series.





# How is DTW Calculated? II

Every possible warping between two time series, is a path though the matrix. We want the best one...

$$DTW(Q,C) = \min\left\{\sqrt{\sum_{k=1}^{K} w_k} \middle| K\right\}$$



This recursive function gives us the minimum cost path

 $\gamma(i,j) = d(q_i,c_j) + \min\{\gamma(i-1,j-1),\gamma(i-1,j),\gamma(i,j-1)\}$ 

# How is DTW Calculated? II

Every possible warping between two time series, is a path though the matrix. We want the best one...

$$DTW(Q,C) = \min\left\{\sqrt{\sum_{k=1}^{K} w_k} \middle| K\right\}$$

This recursive function gives us the minimum cost path

$$\gamma(\mathbf{i},\mathbf{j}) = d(q_{\mathbf{i}},c_{\mathbf{j}}) + \min\{\gamma(\mathbf{i}-1,\mathbf{j}-1),\gamma(\mathbf{i}-1,\mathbf{j}),\gamma(\mathbf{i},\mathbf{j}-1)\}$$



Three Constraints: Given a path W...

- Boundary Constraint: *W* must start and finish in the first and last points of the sequences
- **Continuity**: at any given point in *W*, we can only travel to neighboring points
- Monotonicity: points in *W* must be monotonically ordered



#### Let us visualize the cumulative matrix on a real world problem I





This example shows 2 one-week periods from the power demand time series.

Note that although they both describe 4-day work weeks, the blue sequence had Monday as a holiday, and the red sequence had Wednesday as a holiday.

#### Let us visualize the cumulative matrix on a real world problem II



## Dynamic Time-Warping (Examples)

# So does it work better than Euclidean? Well yes! After all it is more costly..

Dynamic Time Warping



#### Euclidean Distance

18



MIT arrhythmia database



## Results: Error Rate

Dataset	Euclidean	DTW	Using 1-
Word Spotting	4.78	1.10	nearest-
Sign language	28.70	25.93	leaving-
GUN	5.50	1.00	one-out evaluation!
Nuclear Trace	11.00	0.00	
Leaves <sup>#</sup>	33.26	4.07	
(4) Faces	6.25	2.68	
Control Chart*	7.5	0.33	
2-Patterns	1.04	0.00	

## Results: Time (msec)

Dataset	Euclidean	DTW	DTW is
Word Spotting	40	8,600	<sup>215</sup> two to three
Sign language	10	1,110	110 magnitude slower
GUN	60	11,820	197 Euclidean
Nuclear Trace	210	144,470	687 distance
Leaves	150	51,830	345
(4) Faces	50	45,080	901
Control Chart	110	21,900	199
2-Patterns	16,890	545,123	32

#### Fast Approximations to Dynamic Time Warp Distance I



Simple Idea: Approximate the time series with some compressed or downsampled representation, and do DTW on the new representation. How well does this work...



#### Fast Approximations to Dynamic Time Warp Distance II



# Global Constraints

Slightly speed up the calculationsPrevent pathological warpings





A global constraint constrains the indices of the warping path  $w_{\rm k}$  =  $(i,j)_k$  such that  $j-r \le i \le j+r$ 

Where r is a term defining allowed range of warping for a given point in a sequence.

We can learn the best *r* using a training set



Sakoe-Chiba Band

#### Accuracy vs. Width of Warping Window



In general, it's hard to speed up a single DTW calculation



However, if we have to make many DTW calculations (which is almost always the case), we can potentiality speed up the whole process by *lowerbounding*.

Keep in mind that the *lowerbounding* trick works for any situation where you have an expensive calculation that can be *lowerbounded* (string edit distance, graph edit distance etc)



I will explain how *lowerbounding* works in a generic fashion in the next two slides, then show concretely how *lowerbounding* makes dealing with massive time series under DTW possible...

# Lower Bounding I



By definition, for all A, B, we have  $lower\_bound\_distance(A,B) \le DTW(A,B)$
# Lower Bounding II

We can speed up similarity search under DTW by using a lower bounding function

Algorithm Lower_Bound	ing_Sequential_Scan(Q)	Q) (Try to use a cheap
<ol> <li>best_so_far = infini</li> <li>for all sequences in</li> <li>LB_dist = lower_k</li> </ol>	ty; database ound_distance(C <sub>i</sub> , Q);	2); lower bounding calculation as often as possible.
4.       if LB_dist < best_         5.       true_dist = DT         6.       if true_dist < b         7.       best_so_far         8.       best_so_far         9.       endif         10.       endif         11.       endifor	_so_far W(C <sub>i</sub> , Q); est_so_far = true_dist; st_match = i;	Only do the expensive, full calculations when it is absolutely necessary

## Lower Bounding Example



## Lower Bounding Example



## Lower Bounding Example



# Lower Bound of Yi



Yi, B, Jagadish, H & Faloutsos, C. *Efficient retrieval of similar time sequences under time warping*. ICDE 98, pp 23-27.

LB Yi

The sum of the squared length of gray lines represent the minimum corresponding points contribution to the overall DTW distance, and thus can be returned as the lower bounding measure

# Lower Bound of Kim





Kim, S, Park, S, & Chu, W. *An index-based approach for similarity search supporting time warping in large sequence databases.* ICDE 01, pp 607-614 The (maximum) squared difference between the two sequence's first (A), last (D), minimum (B) and maximum points (C) is returned as the lower bound





Envelope-Based Lower Bound

$$LB\_Keogh(Q,C) = \sum_{i=1}^{n} \begin{cases} (C_i - U_i)^2 & \text{if } C_i > U_i \\ (C_i - L_i)^2 & \text{if } C_i < L_i \\ 0 & \text{otherwise} \end{cases}$$



### Lower Bounding the Dynamic Time Warping

#### Minimum Bounding Envelope for bounding the DTW

- Create Minimum Bounding Envelope (MBE) of query Q
- Calculate distance between MBE of Q and any sequence A
- One can show that:  $D(MBE(Q)_{\delta}, A) < DTW(Q, A)$



LB = sqrt(sum([[A > U].\* [A-U]; [A < L].\* [L-A]].^2));

The tightness of the lower bound for each technique is proportional to the length of gray lines used in the illustrations



### Longest Common Subsequence (LCSS)

With Time Warping extreme values (outliers) can destroy the distance estimates. The LCSS model can offer more resilience to noise and impose spatial constraints too.



Everything that is outside the bounding envelope can never be matched

### Longest Common Subsequence (LCSS)

#### LCSS is more resilient to noise than DTW.



#### **Disadvantages of DTW:**

- A. All points are matched
- **B.** Outliers can distort distance
- C. One-to-many mapping

#### Advantages of LCSS:

- A. Outlying values not matched
- B. Distance/Similarity distorted less
- C. Constraints in time & space

## LCSS (Implementation)

Similar dynamic programming solution as DTW, but now we measure similarity not distance.



## Distance Measure Comparison

	AUSLAN	Japanese Vowels	Wafer	ECG	Star Light Curve
Euclidean	.1722	.3838	.08333	.1778	.1767
DTW (full)	.1556	.1063	.0909	.1889	.1647
DTW (window)	.1444	.1063	.0656	.1722	.1566
EDR	.4167	.1417	.3131	.2000	.1606
ERP	.2778	.2970	.0556	.1944	.3936
LCSS	.4185	.1226	.1363	.1278	.1687
LCSS Relaxed	.2259	.1172	.1091	.1278	.1687
TWED	.1741	.1063	.0318	.1278	.1606

Table 5. Error Rates for each Similarity Measure by Dataset

LCSS is the least sensitive to noise because it includes a threshold to define a "match."

# Distance Measure Comparison (Overview)

Method	Complexity	Elastic Matching	One-to-one Matching	Noise Robustness
Euclidean	O(n)	×	$\checkmark$	×
DTW	Ο(n*δ)	$\checkmark$	×	×
LCSS	O(n*δ)	$\checkmark$	$\checkmark$	$\checkmark$





For long time series, shape based similarity will give very poor results. We need to measure similarly based on high level *structure* 



անության հետությանը հայտական է ենչել թնում նունելու չուրելու հետ նրանությունը հետ նրանությունը։ Անհանությունը հետությանը հետությունը հետությանը հետությանը հետությունը հետությունը հետությունը հետությունը հետո



## Structure or Model Based Similarity

The basic idea is to extract global features from the time series, create a feature vector, and use these feature vectors to measure similarity and/ or classify



But which • features? • distance measure/ learning algorithm?



Time Feature Series	A	B	С
Max Value	11	12	19
Autocorrelation	0.2	0.3	0.5
Zero Crossings	98	82	13
ARIMA	0.3	0.4	0.1
• • •	• • •	• • •	• • •

#### Feature-based Classification of Time-series Data

Nanopoulos, Alcock, and Manolopoulos

• features?

• distance measure/ learning algorithm?

#### Learning Algorithm multi-layer perceptron neural network

### Features

mean

variance

skewness

kurtosis

mean (1<sup>st</sup> derivative)

variance (1<sup>st</sup> derivative)

skewness (1st derivative)

kurtosis (1<sup>st</sup> derivative)

#### Learning to Recognize Time Series: Combining ARMA Models with Memory-Based Learning

Deng, Moore and Nechyba

• features?

• distance measure/ learning algorithm?

### Distance Measure

Euclidean distance (between coefficients)

- Used to detect drunk drivers!
- Independently rediscovered and generalized by Kalpakis et. al. and expanded by Xiong and Yeung

#### Features

The parameters of the Box Jenkins model.

More concretely, the coefficients of the ARMA model.

#### Deformable Markov Model Templates for Time Series Pattern Matching Ge and Smyth



### **Compression Based Dissimilarity**

(In general) Li, Chen, Li, Ma, and Vitányi: (For time series) Keogh, Lonardi and Ratanamahatana

- features?
- distance measure/ learning algorithm?

### **Distance Measure** Co-Compressibility $CDM(x, y) = \frac{C(xy)}{C(x) + C(y)}$





### Features

Whatever structure the compression algorithm finds...

The time series is first converted to the SAX symbolic representation\*

### Pattern-Histogram Based Similarity

Lin and Li

• features?

• distance measure/ learning algorithm?

#### **Distance Measure** Euclidean Distance on Bag of Patterns

Clustering using the Bag-of-Patterns approach

8 <mark>(1), 2 (2), </mark>
9 16 reg of the state of the s
19 1
25 0.4 0. 0.0 0. 24 0. 0. 0. 0. 23 0. 0. 0. 0. 0. 0.



#### Features

Represent time series as "bag of patterns" (much like the "bag of words" for text data)

The time series is first converted to the SAX symbolic representation\*

# Summary of Time Series Similarity

- If you have *short* time series, use DTW after searching over the warping window size (and shape), or LCSS
- Then use envelope based lower bounds to speed things up.
- If you have *long* time series, and you know nothing about your data, try compression based dissimilarity or bag-of-patterns similarity.
- If you do know something about your data, try to leverage of this knowledge to extract features.