

COURSE NAME:  
**DATA WAREHOUSING & DATA MINING**

---

# LECTURE 22

## TOPICS TO BE COVERED:

---

- × Time series data
- × Sequence data mining

# MINING TIME-SERIES AND SEQUENCE DATA

## × Time-series database

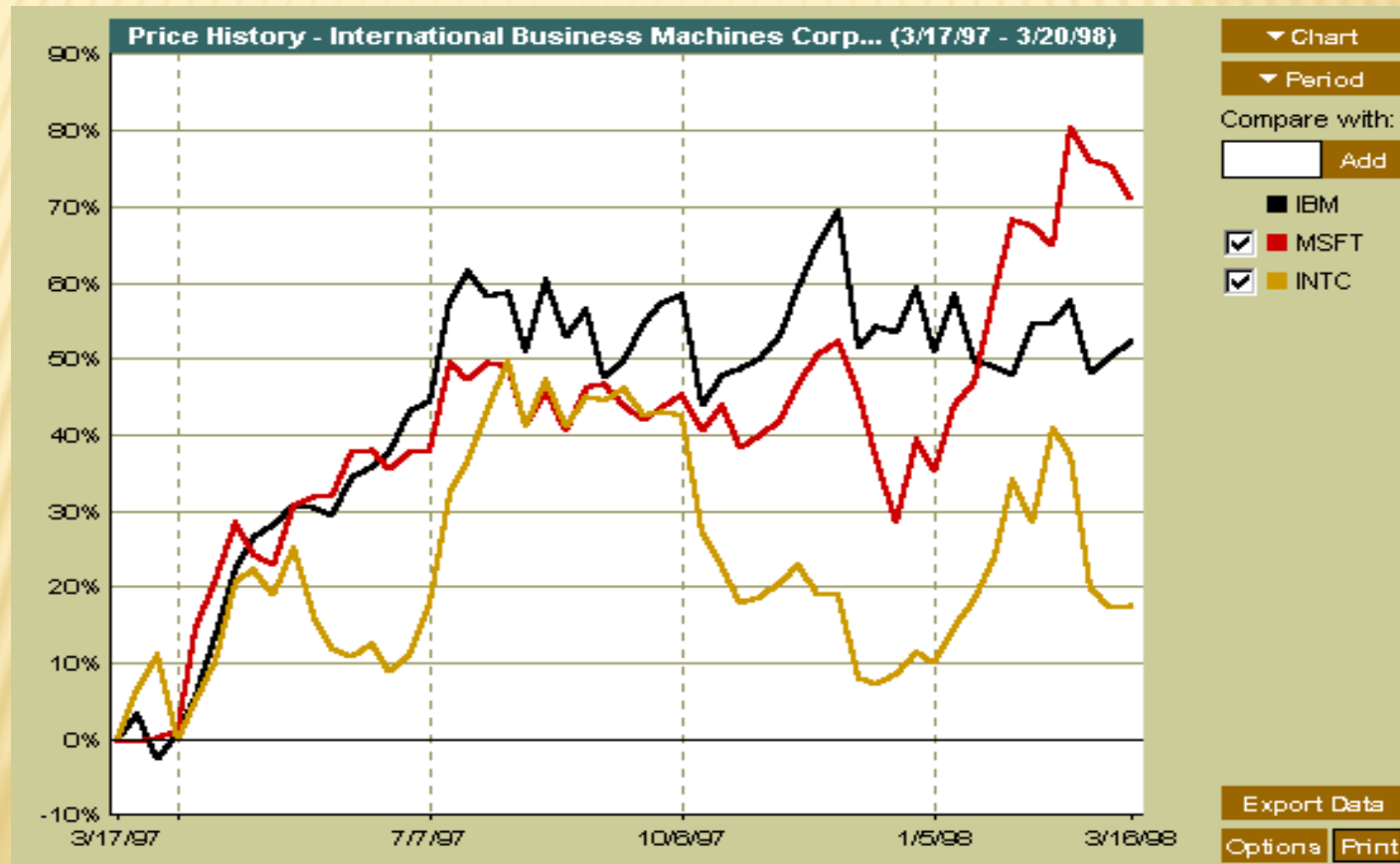
- + Consists of sequences of values or events changing with time
- + Data is recorded at **regular intervals**
- + Characteristic time-series components
  - × Trend, cycle, seasonal, irregular

## × Applications

- + Financial: stock price, inflation
- + Biomedical: blood pressure
- + Meteorological: precipitation

# MINING TIME-SERIES AND SEQUENCE DATA

## Time-series plot



# MINING TIME-SERIES AND SEQUENCE DATA: TREND ANALYSIS

- ✘ A time series can be illustrated as a time-series graph which describes a point moving with the passage of time
- ✘ Categories of Time-Series Movements
  - + Long-term or trend movements (trend curve)
  - + Cyclic movements or cycle variations, e.g., business cycles
  - + Seasonal movements or seasonal variations
    - ✘ i.e, almost identical patterns that a time series appears to follow during corresponding months of successive years.
  - + Irregular or random movements

# ESTIMATION OF TREND CURVE

---

- ✘ The freehand method
  - + Fit the curve by looking at the graph
  - + Costly and barely reliable for large-scaled data mining
- ✘ The least-square method
  - + Find the curve minimizing the sum of the squares of the deviation of points on the curve from the corresponding data points
- ✘ The moving-average method
  - + Eliminate cyclic, seasonal and irregular patterns
  - + Loss of end data
  - + Sensitive to outliers

# DISCOVERY OF TREND IN TIME-SERIES (1)

---

## × Estimation of seasonal variations

### + Seasonal index

- × Set of numbers showing the relative values of a variable during the months of the year
- × E.g., if the sales during October, November, and December are 80%, 120%, and 140% of the average monthly sales for the whole year, respectively, then 80, 120, and 140 are seasonal index numbers for these months

### + Deseasonalized data

- × Data adjusted for seasonal variations
- × E.g., divide the original monthly data by the seasonal index numbers for the corresponding months

# DISCOVERY OF TREND IN TIME-SERIES (2)

---

- ✘ Estimation of cyclic variations
  - + If (approximate) periodicity of cycles occurs, cyclic index can be constructed in much the same manner as seasonal indexes
- ✘ Estimation of irregular variations
  - + By adjusting the data for trend, seasonal and cyclic variations
- ✘ With the systematic analysis of the trend, cyclic, seasonal, and irregular components, it is possible to make long- or short-term predictions with reasonable quality



# SIMILARITY SEARCH IN TIME-SERIES ANALYSIS

---

- ✗ Normal database query finds exact match
- ✗ Similarity search finds data sequences that differ only slightly from the given query sequence
- ✗ Two categories of similarity queries
  - + **Whole Sequence Matching** : find a sequence that is similar to the query sequence
  - + **Subsequence matching**: find all pairs of similar sequences
- ✗ Typical Applications
  - + Financial market
  - + Market basket data analysis
  - + Scientific databases
  - + Medical diagnosis

# ENHANCED SIMILARITY SEARCH METHODS

---

- ✘ Allow for gaps within a sequence or differences in offsets or amplitudes
- ✘ Normalize sequences with amplitude scaling and offset translation
- ✘ Two subsequences are considered similar if one lies within an envelope of  $\varepsilon$  width around the other, ignoring outliers
- ✘ Two sequences are said to be similar if they have enough non-overlapping time-ordered pairs of similar subsequences
- ✘ Parameters specified by a user or expert: sliding window size, width of an envelope for similarity, maximum gap, and matching fraction

# QUERY LANGUAGES FOR TIME SEQUENCES

---

- ✘ Time-sequence query language

- + Should be able to specify sophisticated queries like

Find all of the sequences that are similar to some sequence in class *A*, but not similar to any sequence in class *B*

- + Should be able to support various kinds of queries: range queries, all-pair queries, and nearest neighbor queries

- ✘ Shape definition language

- + Allows users to define and query the overall shape of time sequences

- + Uses human readable series of sequence transitions or macros

- + Ignores the specific details

- ✘ E.g., the pattern **up, Up, UP** can be used to describe increasing degrees of rising slopes

- ✘ Macros: **spike, valley**, etc.

# SEQUENTIAL PATTERN MINING

---

- ✘ Mining of frequently occurring patterns related to time or other sequences
- ✘ Sequential pattern mining usually concentrate on symbolic patterns
- ✘ Examples
  - + Renting “Star Wars”, then “Empire Strikes Back”, then “Return of the Jedi” in that order
  - + Collection of ordered events within an interval
- ✘ Applications
  - + Targeted marketing
  - + Customer retention
  - + Weather prediction

# MINING SEQUENCES (CONT.)

## *Customer-sequence*

CustId	Video sequence
1	{(C), (H)}
2	{(AB), (C), (DFG)}
3	{(CEG)}
4	{(C), (DG), (H)}
5	{(H)}

## *Map Large Itemsets*

Large Itemsets	MappedID
(C)	1
(D)	2
(G)	3
(DG)	4
(H)	5

*Sequential patterns with support > 0.25*

*{(C), (H)}*

*{(C), (DG)}*

# PERIODICITY ANALYSIS

---

- × Periodicity is everywhere: tides, seasons, daily power consumption, etc.
- × **Full periodicity**
  - + Every point in time contributes (precisely or approximately) to the periodicity
- × **Partial periodicity**: A more general notion
  - + Only some segments contribute to the periodicity
    - × Jim reads NY Times 7:00-7:30 am every week day
- × **Cyclic association rules**
  - + Associations which form cycles
- × **Methods**
  - + Full periodicity: FFT, other statistical analysis methods
  - + Partial and cyclic periodicity: Variations of Apriori-like mining methods