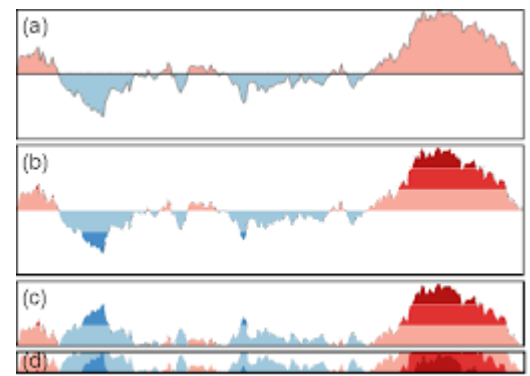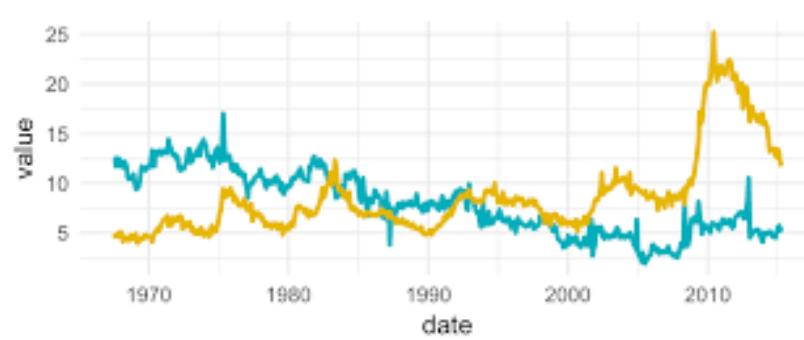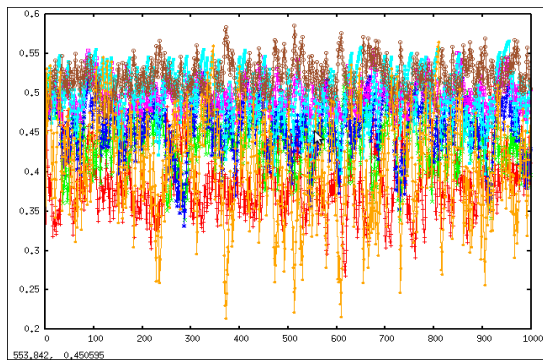# Mining Time Series

Mining Massive Datasets
Prof. Carlos Castillo
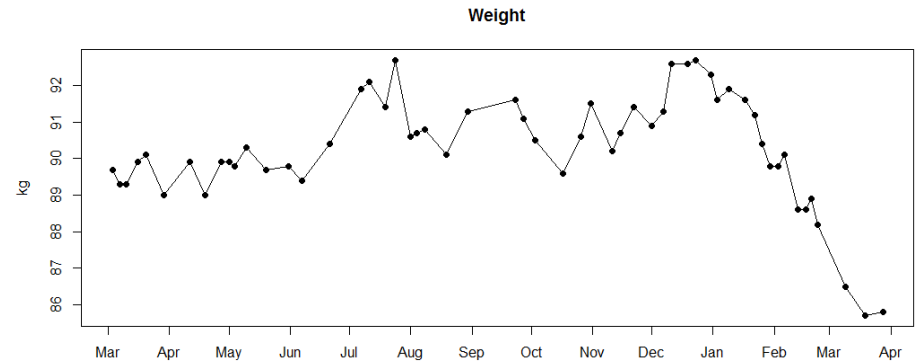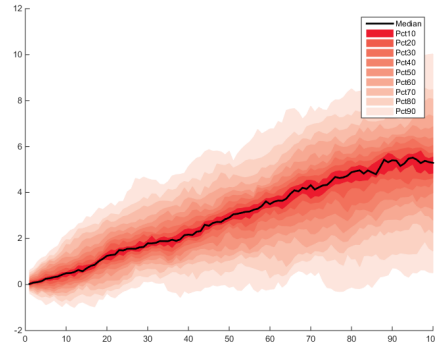Topic 27

# Sources

- Data Mining, The Textbook (2015) by Charu Aggarwal (chapter 14)

- Introduction to Time Series Mining (2006) tutorial by Keogh Eamonn [alt. link]

- Time Series Data Mining (2006) slides by Hung Son Nguyen

# Why do we mine time series? Examples

# Seismic data

- Observations = earthquakes

- Goal: characterize when peeks occur





Earthquake sequence in time

Magnitude

# Liquid metal droplets

◇ = length of hot metal droplet

■ = droplet release (chaotic, noisy)

Goal: prediction of release

5

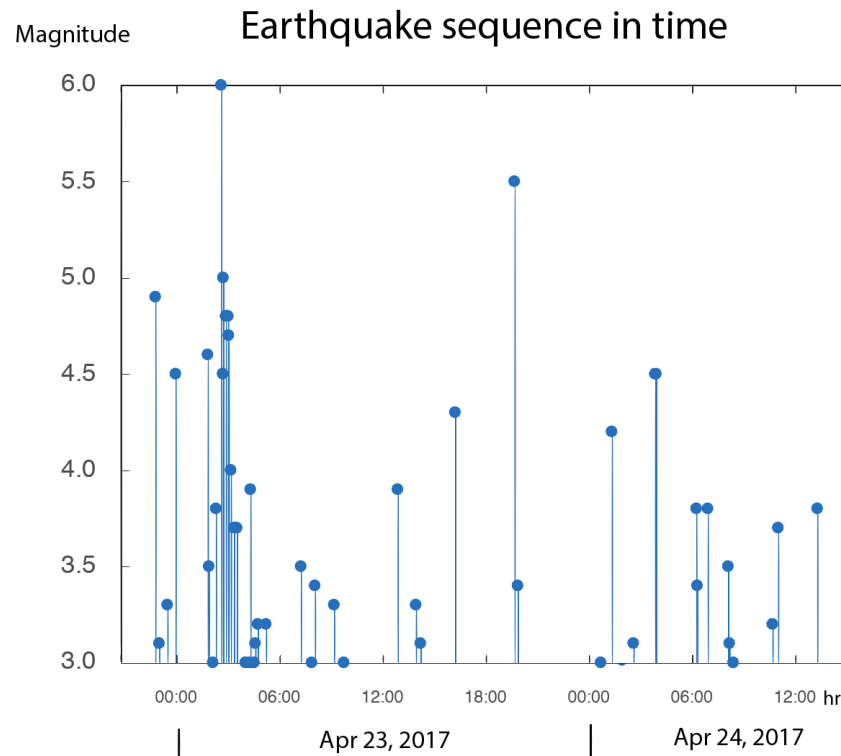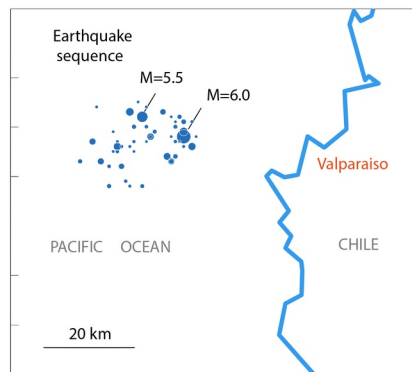# Stock prices

Price →



Volume traded →

Goal: find hidden patterns providing an advantage

6

# Video data / gestures

- Series of angles of articulations in the body

- Temporal patterns can reveal gestures



Point

Steady pointing

Hand moving to shoulder level

Hand at rest

Gun-Draw

Steady pointing

Hand moving to shoulder level

Hand moving down to grasp gun

Hand moving above holster

Hand at rest

# Applications

- Clustering

- Classification

- Motif discovery

- Event detection

- ...

1) All require a reasonable definition of the **similarity** between two time series

2) All can be done in **real-time** or **retrospectively**

# Context vs Behavior

- **Contextual attribute(s)**
  - $x(i)$ = $t_i$ = timestamp is the typical one
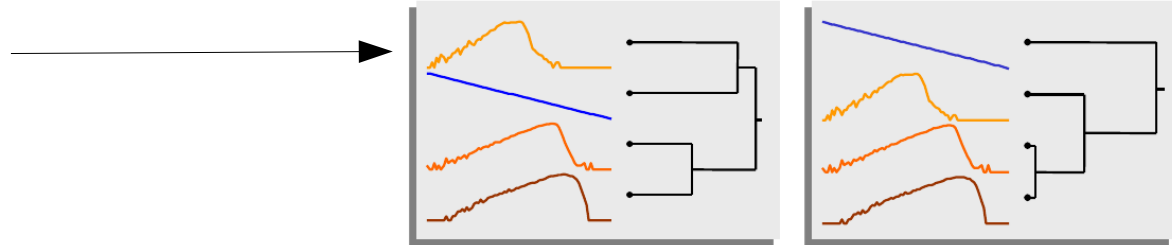  - Sometimes other attributes providing context
- **Behavioral attribute(s)**
  - $y^j(i)$ = temperature, angle, price, sensor reading, …
    $j \in 1 \ldots d$

# What are the difficulties?

- High sampling rate of many series over extended periods of time means ...

  – Tons of data

  – Things are bound to fail at several points (missing data, noisy data)

- Subjectivity

# Preparing a time series
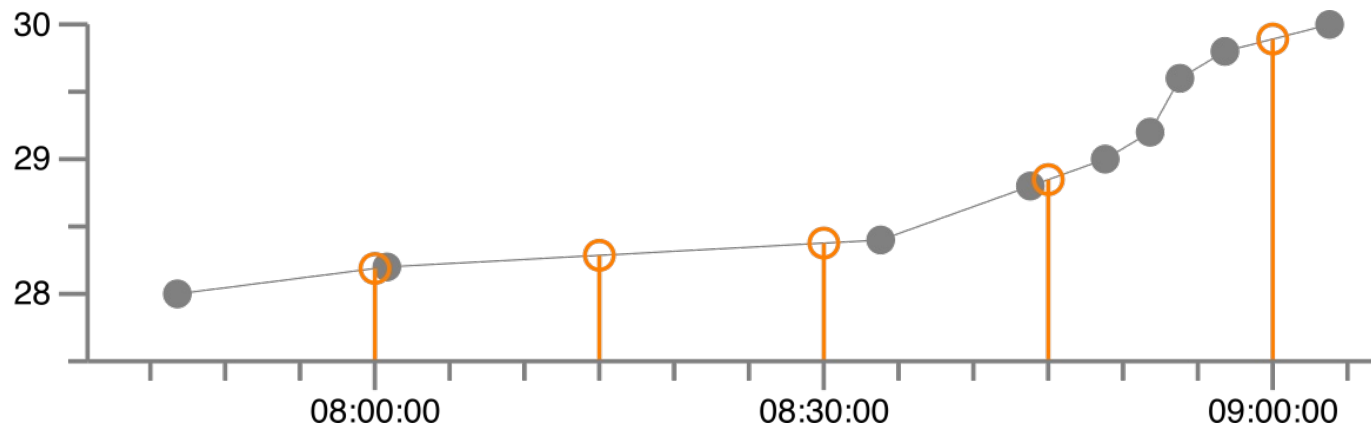
# Notation: multivariate time series

- Length $n$, timestamps $t_1, t_2, \ldots, t_n$

- Values at time $t_i : (y_i^1, y_i^2, \ldots, y_i^d)$

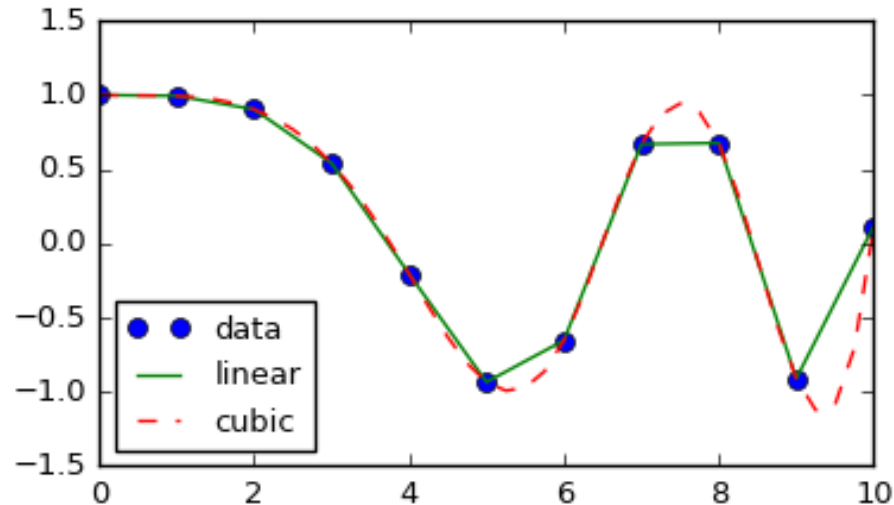- If series is univariate we drop the superscript

# **Missing values**: linear interpolation

- Let $t_i < t_x < t_j$ $\qquad y_x = y_i + \left( \dfrac{t_x - t_i}{t_j - t_i} \right) \cdot (y_j - y_i)$

- Example: make an irregular series regular

# **Missing values**: splines
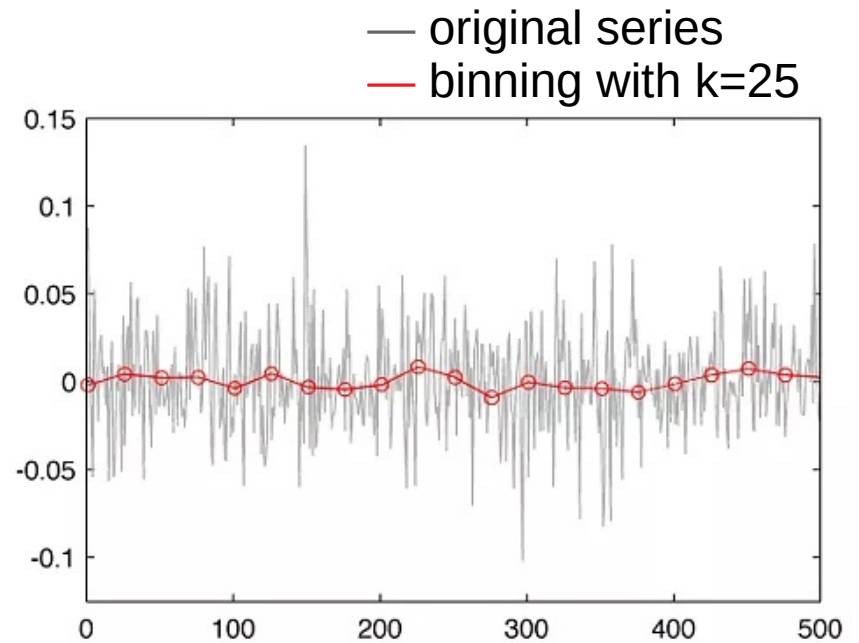
Cubic polynomials between $y_i$, $y_{i+1}$ that have the same slope at those points as the original curve.

# **Noise removal**: binning

- Replace series by average of values in bins (subsequences) of length k

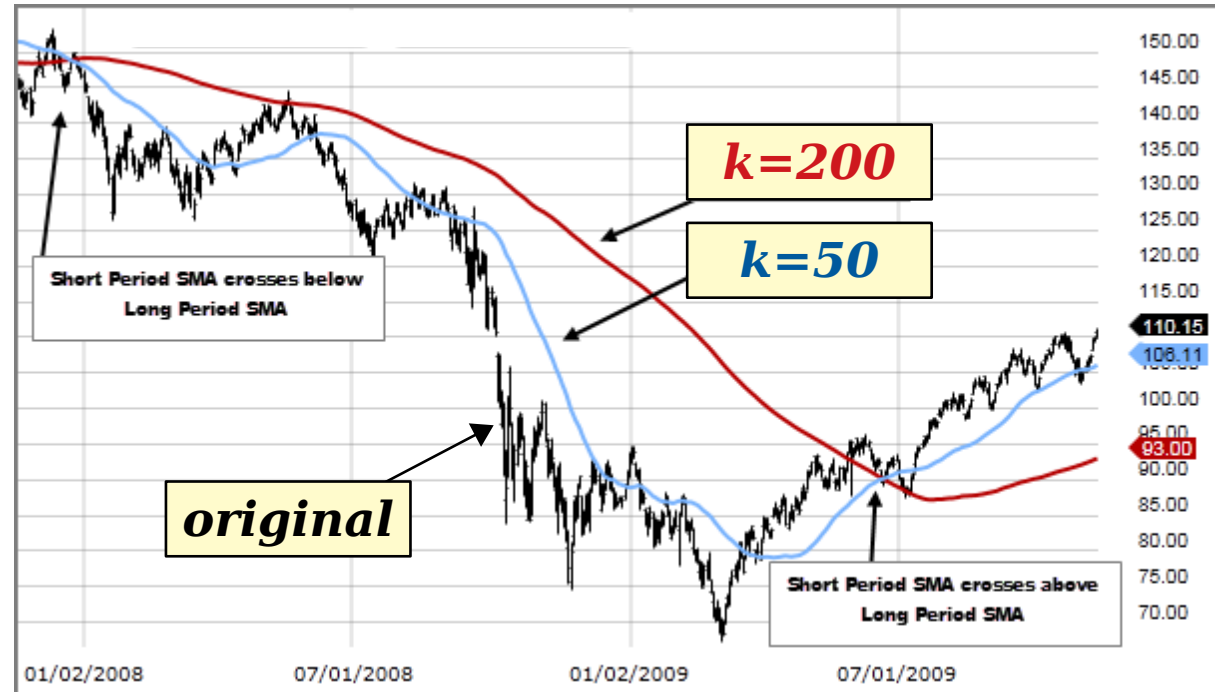$$y'_{i+1} = \frac{1}{k} \sum_{r=1}^{k} y_{i \cdot k + r}$$



— original series
— binning with k=25

http://www.quantatrisk.com/2013/03/22/rebinning-of-financial-time-series/

# Noise removal:
## moving average smoothing

- Equivalent to overlapping bins

$$y'_i = \frac{1}{k} \sum_{r=1}^{k} y_{i-r+1}$$

- Larger k leads to smoother series, but losses more information

- Use smaller k for first k-1 items

**k=200**

**k=50**

Short Period SMA crosses below Long Period SMA

*original*

Short Period SMA crosses above Long Period SMA

110.15
108.11
93.00

150.00
145.00
140.00
135.00
130.00
125.00
120.00
115.00
100.00
95.00
90.00
85.00
80.00
75.00
70.00

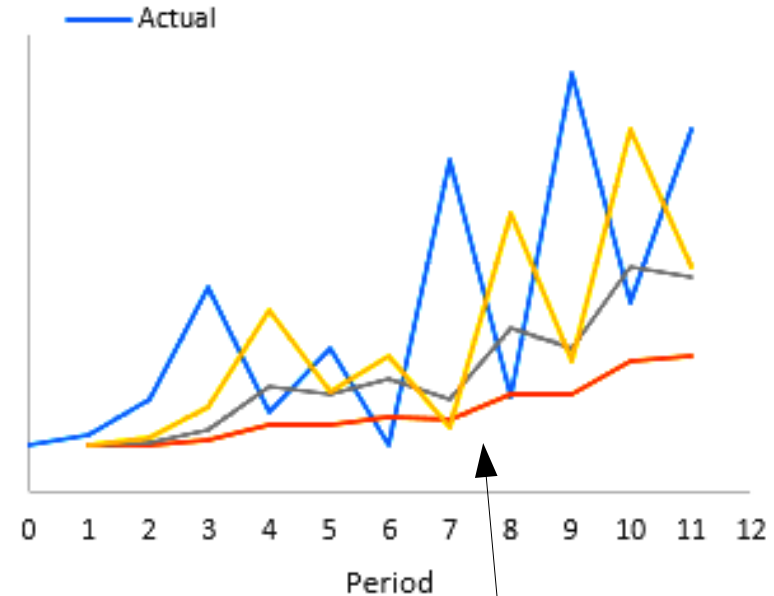01/02/2008     07/01/2008     01/02/2009     07/01/2009

16

# **Noise removal**: exponential smoothing

- Combine previously smoothed point with current point

$$y_i' = \alpha \cdot y_i + (1 - \alpha) \cdot y_{i-1}'$$

- Recursively substituting

$$y_i' = (1 - \alpha)^i \cdot y_0' + \alpha \sum_{j=1}^{i} y_j \cdot (1 - \alpha)^{i-j}$$
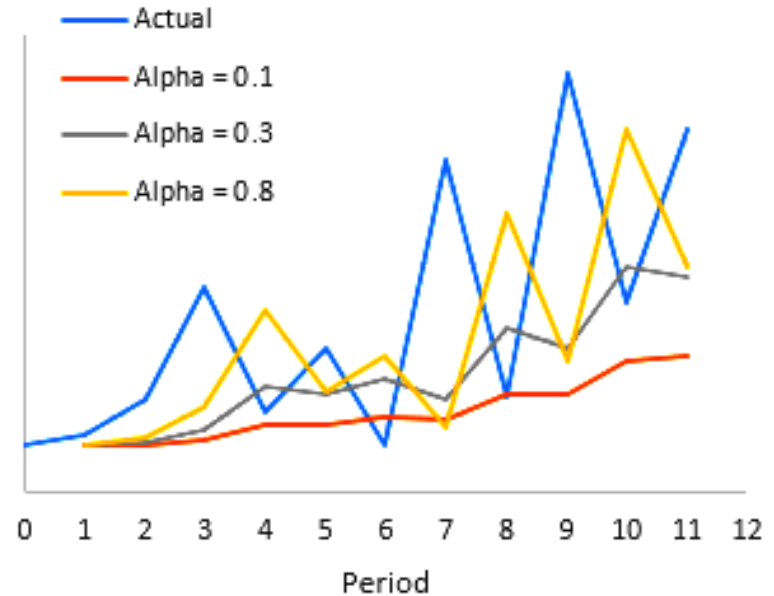
Which y' has the larger alpha?

# **Noise removal**: exponential smoothing

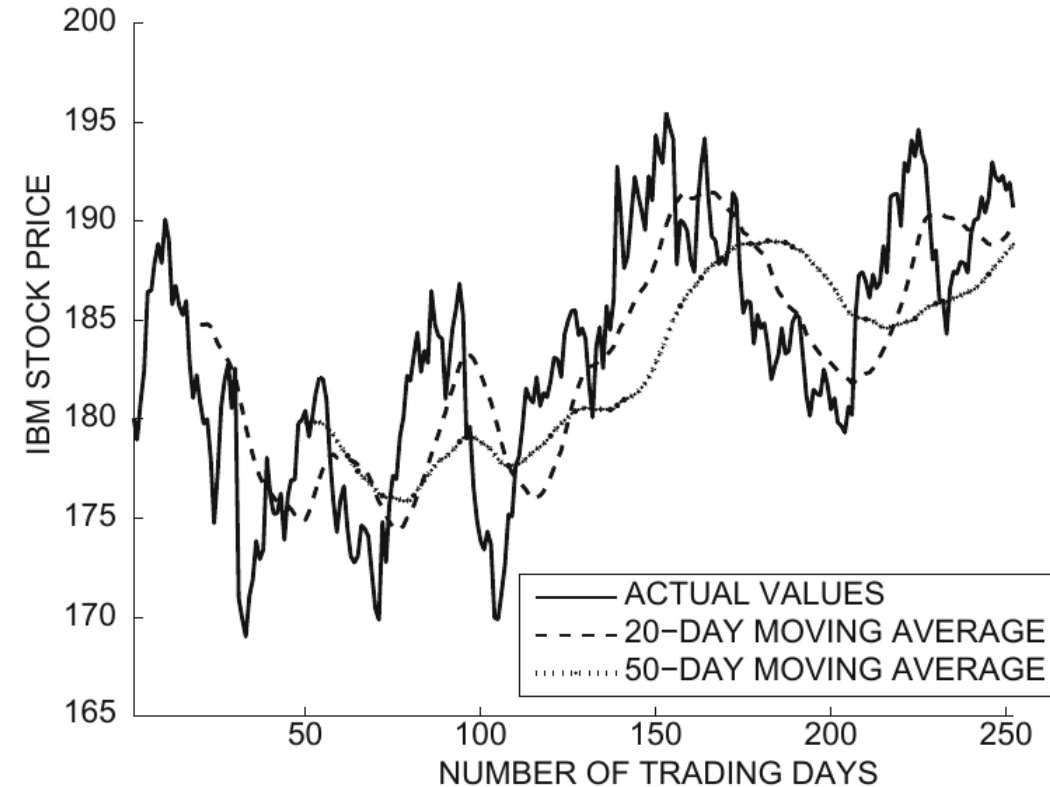- Combine previously smoothed point with current point

$$y_i' = \alpha \cdot y_i + (1 - \alpha) \cdot y_{i-1}'$$
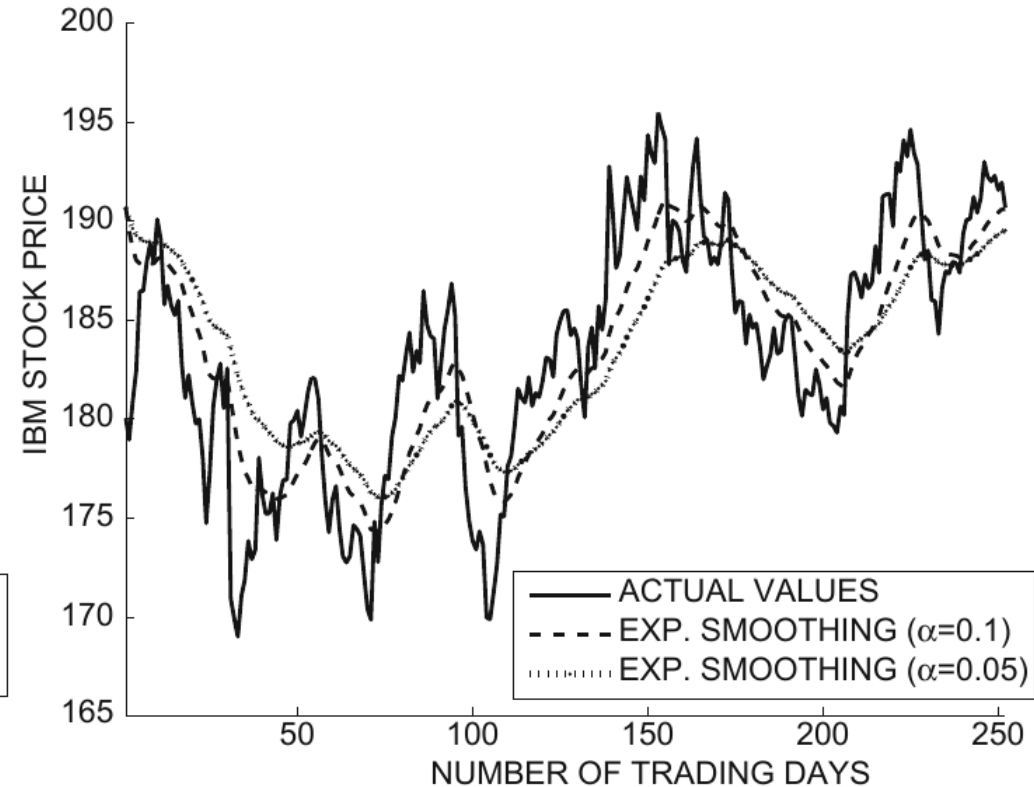
- Recursively substituting

$$y_i' = (1 - \alpha)^i \cdot y_0' + \alpha \sum_{j=1}^{i} y_j \cdot (1 - \alpha)^{i-j}$$



Actual
Alpha = 0.1
Alpha = 0.3
Alpha = 0.8

0   1   2   3   4   5   6   7   8   9   10   11   12

Period

# Moving average vs exponential smoothing



(a) Moving average smoothing

(b) Exponential smoothing

# Exercise

- Given the following series:

| t | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| y(t) | 2 | 4 | 12 | 2 | 1 | -2 | 0 | 15 | 3 | 3 |
| 1. y'(t) | | | | | | | | | | |
| 2. y'(t) | | | | | | | | | | |

- 1. Moving average with k=3
- 2. Exponential average with alpha=0.5

20

# Summary

# Things to remember

- Series preparation
  - Interpolation
  - Smoothing

# Exercises for TT27-TT29

- Data Mining, The Textbook (2015) by Charu Aggarwal
    - Exercises 14.10 → 1-6