# Mining Time Series

## Mining Massive Datasets

Materials provided by Prof. Carlos Castillo — https://chato.cl/teach

Instructor: Dr. Teodora Sandra Buda — https://tbuda.github.io/
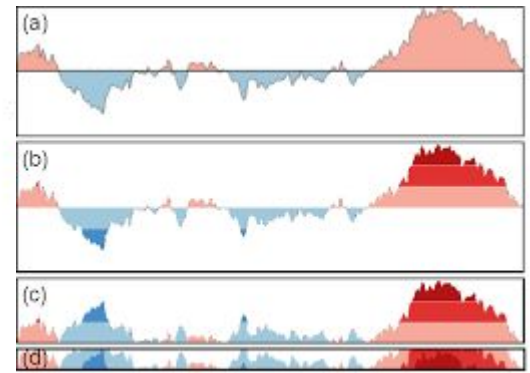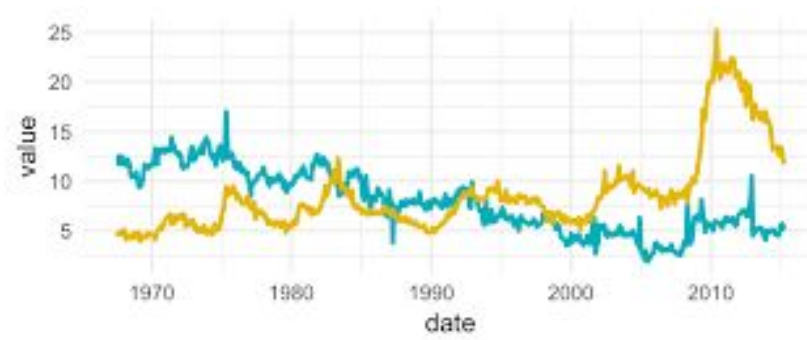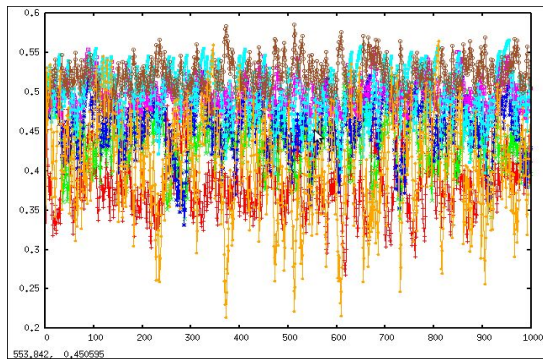
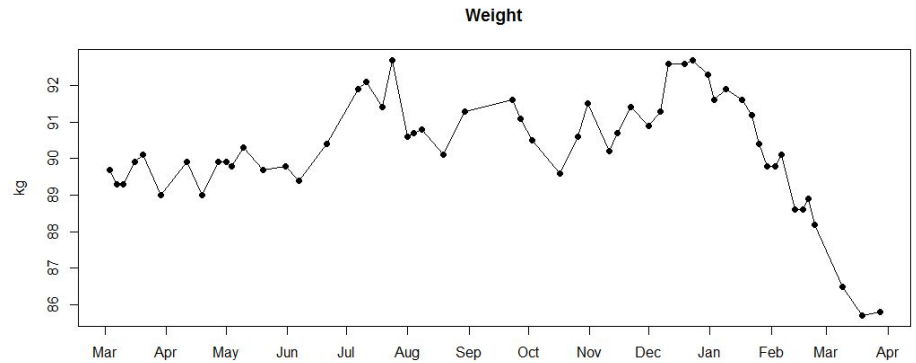IF YOUR DATA HAS A TIME STAMP
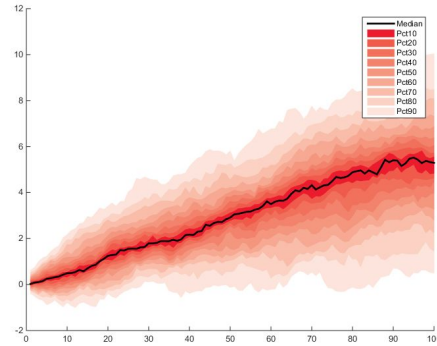
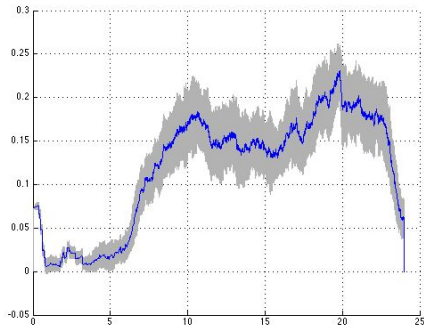YOU'RE A TIME SERIES ANALYST, HARRY

memegenerator.net

# 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
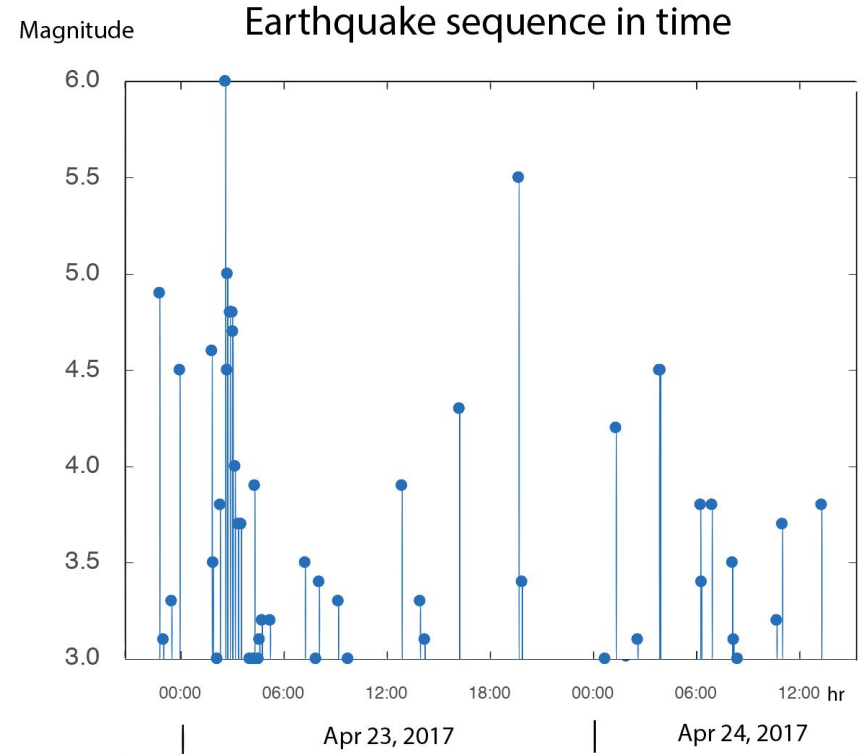
# Seismic data

- Observations = earthquakes

- Goal: characterize when peaks occur

# Liquid metal droplets

◇ = length of hot
metal droplet

■ = droplet release
  – (chaotic, noisy)

Goal: prediction of
release

# Stock prices



Price

Volume traded

Goal: find hidden
patterns providing an
advantage

**BEYOND MEAT (BYND)** STOCK NAS

▲ 81.72 USD 5.96 (7.97%) 02:41:57 PM EDT BTT

# Video data / gestures

- Series of angles of articulations in the body

- Temporal patterns can reveal gestures

# 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, …

  - $j \in 1 \dots d$

# 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 $t_1$, $t_2$, ..., $t_n$
- Values at time $t_i$ : $(y_i^1, y_i^2, ..., y_i^d)$
- If series is univariate we drop the superscript

# **Missing values**: linear interpolation

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

- 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
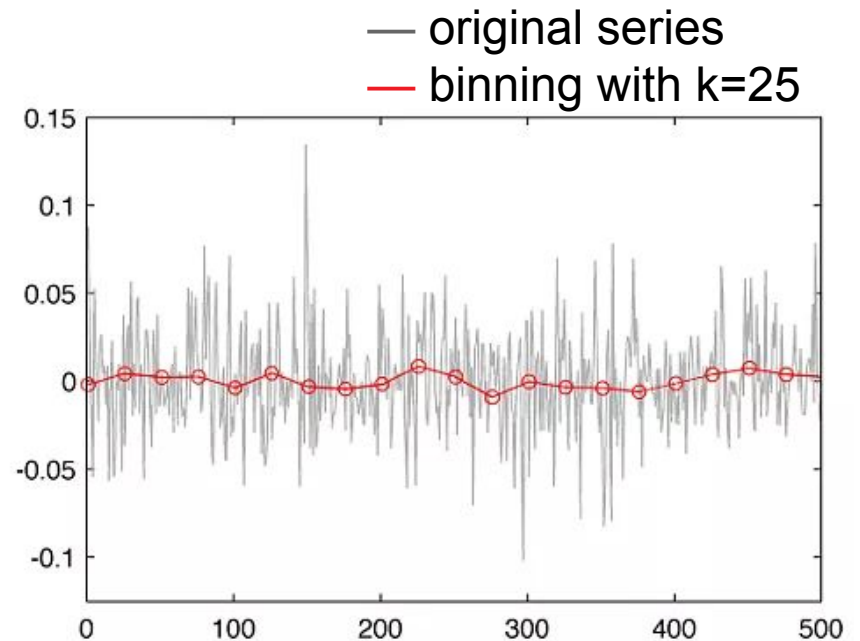
k=200

k=50

Short Period SMA crosses below Long Period SMA

*original*

Short Period SMA crosses above Long Period SMA

150.00
145.00
140.00
135.00
130.00
125.00
120.00
115.00
110.15
108.11
100.00
95.00
93.00
90.00
85.00
80.00
75.00
70.00

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

# **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



Period
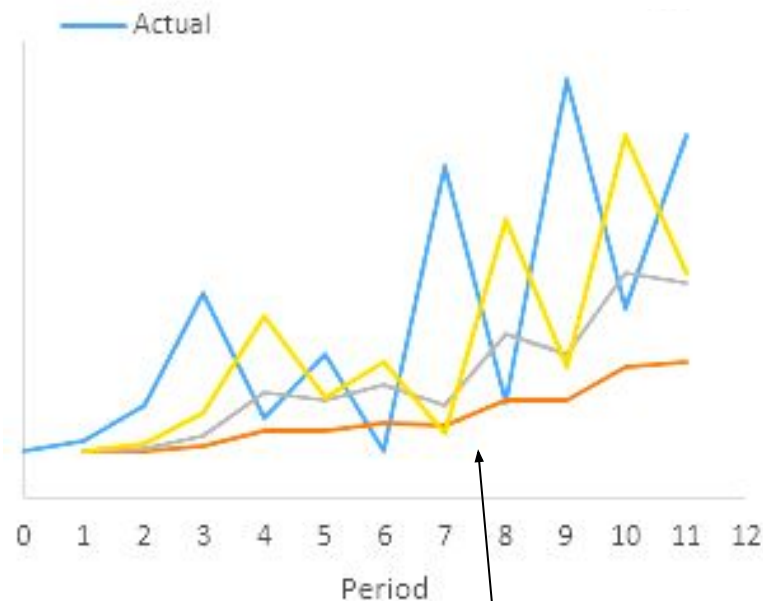
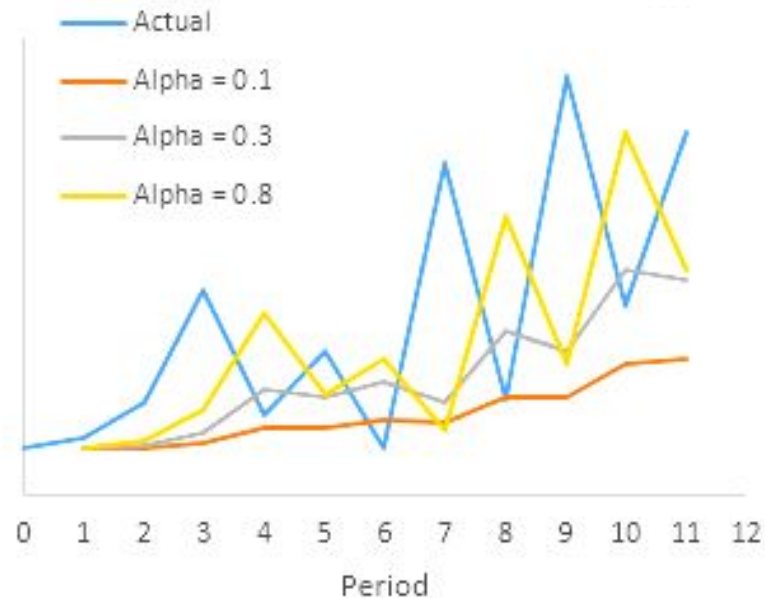**Which y' has the larger alpha?**

# **Noise removal**: exponential smoothing

- Combine previously smoothed point with
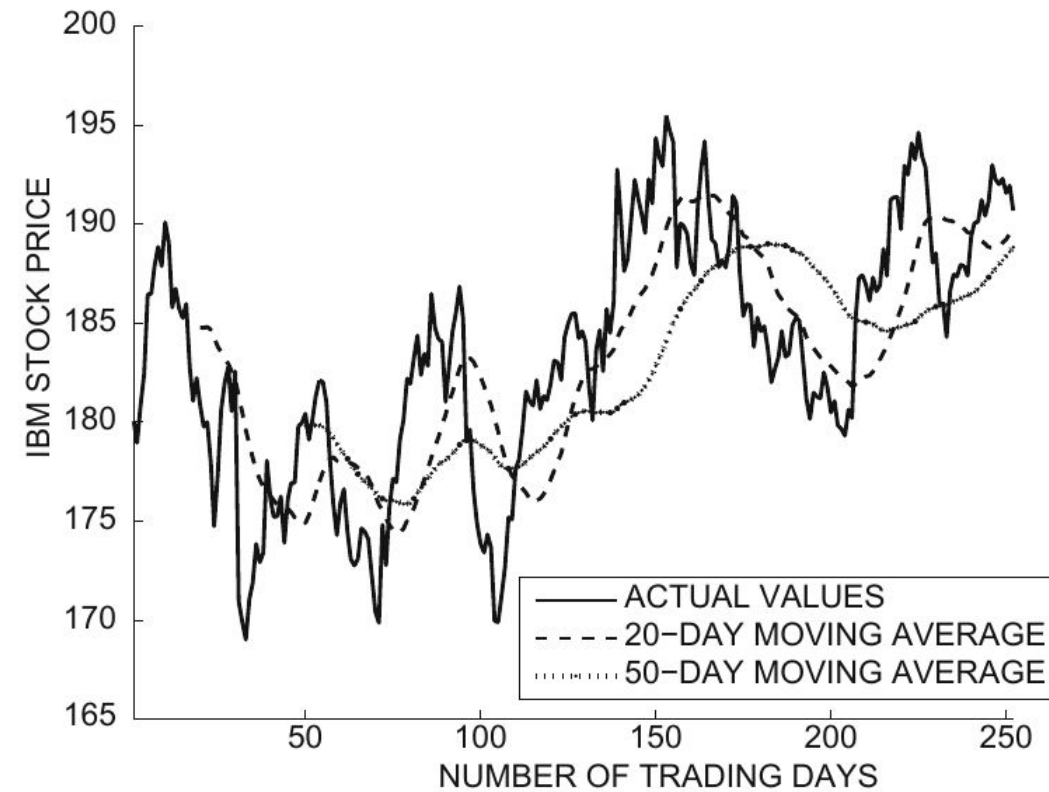  current point

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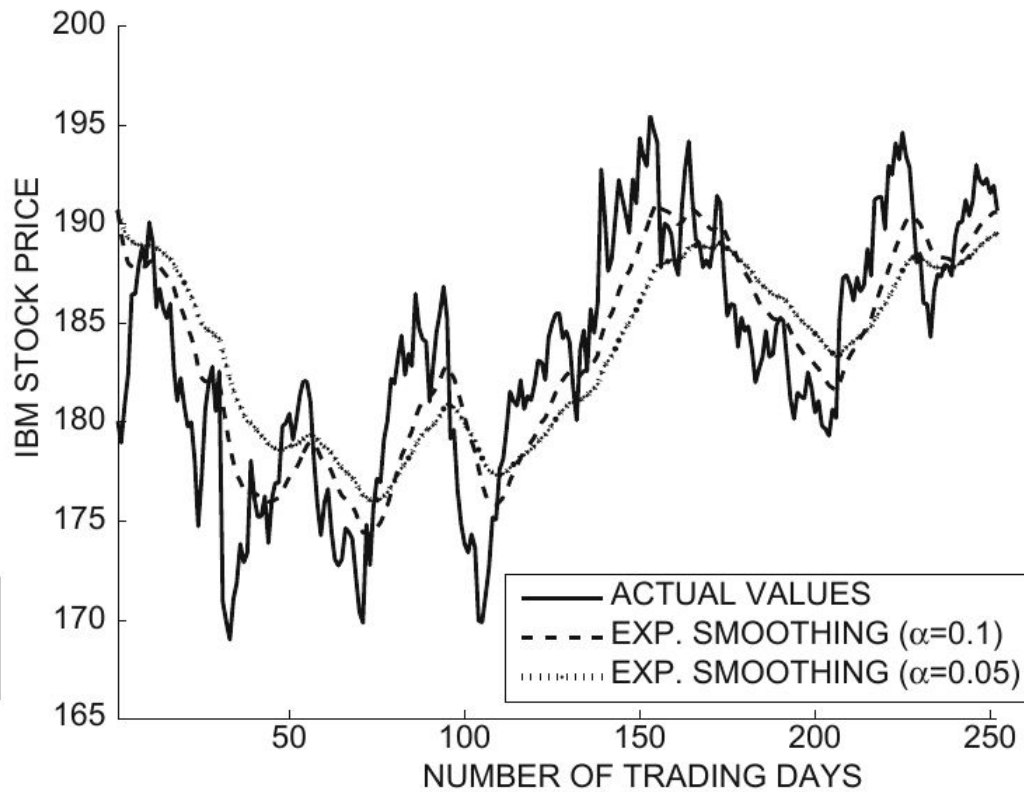
- Recursively substituting

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

# Moving average vs exponential smoothing



(a) Moving average smoothing

(b) Exponential smoothing

# Exercise: smooth a time series

- 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

# Answer

- 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 |
| $y_t'$ | 2 | 3 | 6 | 6 | 5 | 0.33 | -0.33 | 4.33 | 6 | 7 |
| $y_t''$ | 2 | 3 | 7.5 | 4.75 | 2.88 | 0.44 | 0.22 | 7.61 | 5.30 | 4.15 |

- $y_t'$: moving average with k=3
- $y_t''$: exponential average with alpha=0.5

# Answer (code)

python

```python
x = [2, 4, 12, 2, 1, -2, 0, 15, 3, 3]
```

```python
k = 3
y = [0] * len(x)
for i in range(len(x)):
    s = 0
    c = 0
    for j in range(k):
        if i-j >= 0:
            s = s + x[i-j]
            c += 1
    y[i] = s / c if c > 0 else 0
```

# 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