

Curs 11-12:

Analiza seriilor de timp

# Structura

- Motivație
- Pre-procesarea seriilor de timp
- Predicție
- Identificare șabloane
- Detecție anomalii

# Motivație

**Problema:** Se cunosc date săptămânale privind indexul Dow Jones și se dorește identificarea acțiunilor pentru care creșterea de profit va fi cea mai mare în săptămâna care urmează

**Set date:** Dow Jones Index (UCI Machine Learning - Brown, Pelosi & Dirksa, 2013) - 750 înregistrări, 16 attribute

Exemple de companii cotate și pt care sunt înregistrate informații:

3M	MMM	Cisco Systems	CSCO
American Express	AXP	Coca-Cola	KO
Alcoa	AA	DuPont	DD
AT&T	T	ExxonMobil	XOM
Bank of America	BAC	General Electric	GE
Boeing	BA	Hewlett-Packard	HPQ
Caterpillar	CAT	The Home Depot	HD
Chevron	CVX	Intel	INTC
			IBM

# Motivație

**Întrebare:** care acțiune va înregistra cea mai mare creștere în săptămâna care urmează?

Exemplu set de date [Dow Jones Index from <http://archive.ics.uci.edu/ml/datasets.html>]  
16 attribute

- **quarter:** the yearly quarter (1 = Jan-Mar; 2 = Apr-Jun; 3 = Jul-Sep; 4 = Oct-Dec).
- **stock:** the stock symbol (lista de pe slide-ul anterior)
- **date:** the last business day of the work (de obicei e Vineri)
- **open:** the price of the stock at the beginning of the week
- **high:** the highest price of the stock during the week
- **low:** the lowest price of the stock during the week
- **close:** the price of the stock at the end of the week
- **volume:** the number of shares of stock that traded hands in the week
- **percent\_change\_price:** the percentage change in price throughout the week
- **percent\_change\_volume\_over\_last\_wek:** the percentage change in the number of shares of stock that traded hands for this week compared to the previous week
- **previous\_weeks\_volume:** the number of shares of stock that traded hands in the previous week

# Motivație

Întrebare: care acțiune va înregistra cea mai mare creștere în săptămâna care urmează?

Exemplu [Dow Jones Index from <http://archive.ics.uci.edu/ml/datasets.html>]  
16 attribute

- **next\_weeks\_open**: the opening price of the stock in the following week
- **next\_weeks\_close**: the closing price of the stock in the following week
- **percent\_change\_next\_weeks\_price**: the percentage change in price of the stock in the following week
- **days\_to\_next\_dividend**: the number of days until the next dividend
- **percent\_return\_next\_dividend**: the percentage of return on the next dividend

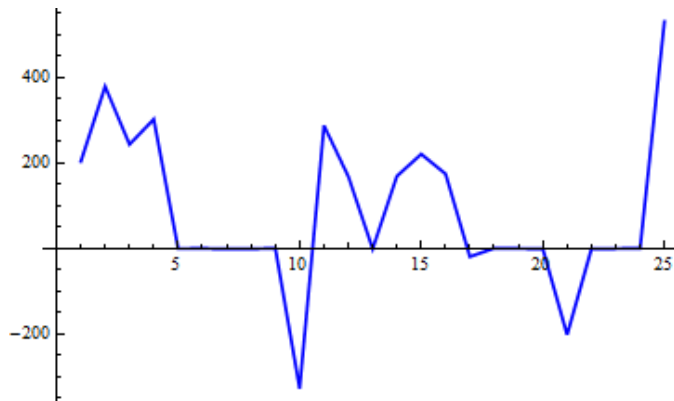
# Motivație

**Întrebare:** care acțiune va înregistra cea mai mare creștere în săptămâna care urmează?

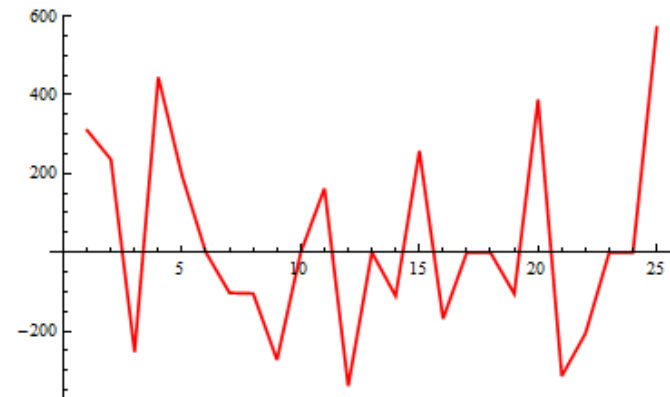
Exemplu [Dow Jones Index de la <http://archive.ics.uci.edu/ml/datasets.html>]

16 attribute

**percent\_change\_next\_weeks\_price:** the percentage change in price of the stock in the following week



IBM



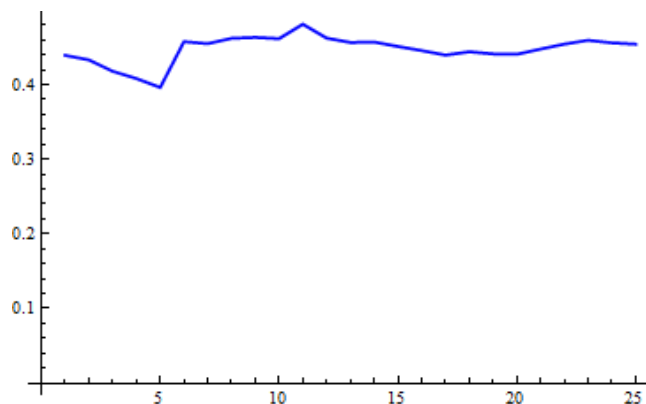
HP

# Motivație

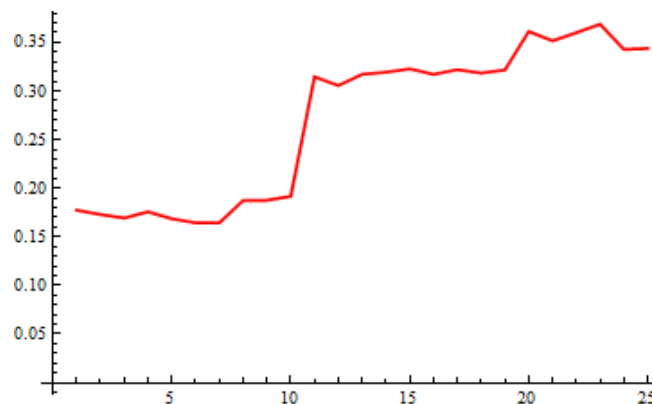
**Întrebare:** care acțiune va înregistra cea mai mare creștere în săptămâna care urmează?

Exemplu [Dow Jones Index de la <http://archive.ics.uci.edu/ml/datasets.html>]  
16 attribute

**percent\_return\_next\_dividend:** the percentage of return on the next dividend



IBM



HP

# Motivație

Pe lângă datele financiare există o mulțime de alte surse de serii de timp:

- **Senzori:**

- Date de mediu colectate prin intermediul diferitelor tipuri de senzori (ex: temperatura, presiune, umiditate)
- Date colectate de la dispozitive de producție sau ustensile de uz casnic (ex: consum de energie)

- **Date medicale**

- Electrocardiograma (ECG)
- Electroencefalograma (EEG)
- Date de monitorizare în timp real a pacienților de la terapie intensivă

- **Date de tip “web log” (clickstream data)**

- Secvențe indicând vizite ale unor pagini web



# Motivație

Pe lângă datele financiare există o mulțime de alte surse de serii de timp:

- **Senzori:**

- Date de mediu colectate prin intermediul diferitelor tipuri de senzori (ex: temperatura, presiune, umiditate)
- Date colectate de la dispozitive de producție sau ustensile de uz casnic (ex: consum de energie)

Task: **predicție valori viitoare**

- **Date medicale**

- Electrocardiograma (ECG)
- Electroencefalograma (EEG)
- Date de monitorizare în timp reali a pacienților de la terapie intensivă

Task: **identificare comportament anormal**

- **Date de tip “web log” (clickstream data)**

- Secvențe indicând vizite ale unor pagini web

Task: **identificare tipare de utilizare, profile de utilizatori**

# Serii de timp

**Exemplu 1** (percentage of return on the next dividend for first 10 weeks included in Dow Jones Index dataset)

0.177, 0.172, 0.169, 0.175, 0.168, 0.164, 0.164, 0.187, 0.187, 0.191

**Obs:** Momentul de timp nu apare ca variabilă explicită. Totuși valorile specificate trebuie interpretate în contextul unor momente de timp.

- Timpul este atribut **contextual**
- Valoarea înregistrată este atribut **comportamental**

**Exemplu 2** (temperatura la prânz înregistrată în 7 zile consecutive)

21, 24, 23, 25, 22, 19, 20

Atributul contextual este timpul, cel comportamental este temperatura

# Tipuri de serii

Există diferite tipuri de serii de timp (temporale)

In raport cu domeniul de timp:

- Continue (e.g. EEG)
- Discrete (denumite secvențe)

In raport cu attributele comportamentale

- Univariate (un atribut)
- Multivariate sau vectoriale (mai multe attribute)

# Tipuri de prelucrări

- Pre-procesări

- Tratare valori absente
- Eliminare zgomot

- Analiză

- Identificare tendința
- Identificare caracter sezonier (comportament care se manifestă periodic în anumite intervale de timp)

- Predicție (prognoză)

- Estimarea valorilor corespunzătoare unor momente ulterioare de timp

- Detecție șabloane și anomalii

- Identificarea unor fragmente din serie care respectă un tipar cunoscut sau un tipar identificat în timpul analizei
- Identificarea unor serii care au comportament semnificativ diferit

# Pre-procesarea seriilor de timp

## Valori absente

### Problema:

- Lipsesc valori corespunzătoare unor momente de timp (de exemplu din cauza unor defecte ale senzorilor)
- In special când sunt mai multe attribute comportamentale (colectate de la senzori independenți) ar trebui asigurată sincronizarea între serii prin completarea valorilor absente

### Soluție:

- Valoarea absentă este estimată folosind interpolare
- Caz simplu: [interpolare liniară](#)

# Pre-procesarea seriilor de timp

## Imputarea valorilor absente prin interpolare liniară

Fie  $(y_1, y_2, \dots, y_n)$  o serie de timp corespunzătoare momentelor  $(t_1, t_2, \dots, t_n)$

- Presupunem că lipsește valoarea corespunzătoare momentului  $t_{i+1}$ .
- Presupunând în plus că atributul comportamental  $y$  variază liniar în raport cu  $t$  pe intervalul  $[t_i, t_{i+2}]$  se poate estima valoarea lui  $y_{i+1}$

$$y_{i+1} = y_i + \frac{t - t_i}{t_{i+2} - t_i} (y_{i+2} - y_i)$$

# Pre-procesarea seriilor de timp

## Eliminarea/reducerea zgomotului

**Problema:** dispozitivele utilizate pt colectarea datelor (senzorii) pot fi afectate de bruiaje, a.î. seria poate conține valori generate în mod artificial în procesul de colectare a datelor și care nu reflectă comportamentul real al atributului înregistrat

## Metode:

- Compactare (**Binning**)
- Netezire (**Moving-Average Smoothing**)

# Pre-procesarea seriilor de timp

## Binning

### Idee:

- Intervalul de timp global  $[t_1, t_n]$  corespunzător seriei  $(y_1, y_2, \dots, y_n)$  este divizat în  $m$  subintervale conținând fiecare câte  $k$  elemente ( $m=n/k$ )
- Fiecare subinterval va fi asociat unei valori calculate ca medie a valorilor din seria de timp ce corespund momentelor incluse în subinterval

### Observații:

- Se presupune că momentele de timp corespunzătoare seriei inițiale sunt egal distanțate
- Se **reduce numărul de valori disponibile de  $k$  ori** (este un tip de compresie cu pierdere de informație)

$$(t_1, t_2, \dots, t_n) \rightarrow ((t_1, \dots, t_k), (t_{k+1}, \dots, t_{2k}), \dots, (t_{(m-1)k+1}, \dots, t_{mk}))$$

$$(y_1, y_2, \dots, y_n) \rightarrow (z_1, z_2, \dots, z_m)$$

$$z_i = \frac{1}{k} \sum_{j=1}^k y_{(i-1)k+j}, \quad i = \overline{1, m}$$



# Pre-procesarea seriilor de timp

## Moving average smoothing

**Idee:** se reduce pierderea de informație cauzată de binning folosind “ferestre” de mediere care se suprapun, adică media se calculează pentru elementele ce aparțin unei ferestre mobile (se deplasează de-a lungul seriei)

$$(t_1, t_2, \dots, t_n) \rightarrow ((t_1, \dots, t_k), (t_2, \dots, t_{k+1}), \dots, (t_{(m-1)k+1}, \dots, t_{mk}))$$

$$(y_1, y_2, \dots, y_n) \rightarrow (z_1, z_2, \dots, z_{n-k+1})$$

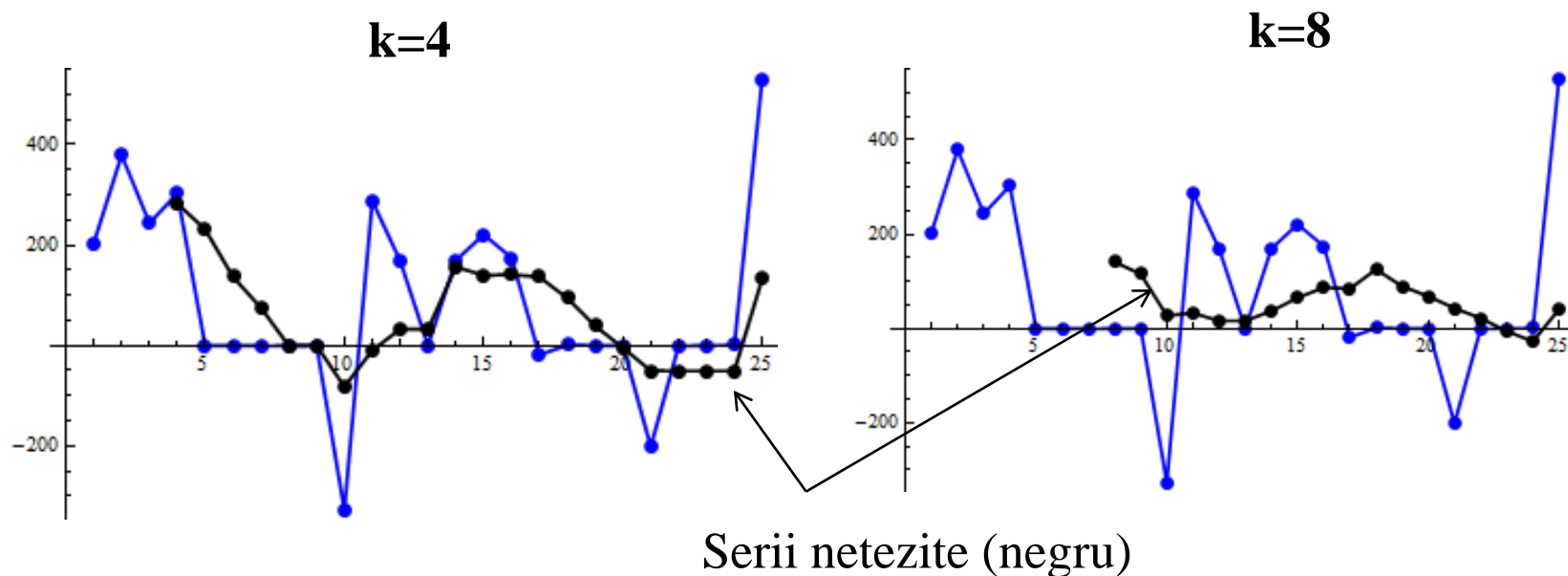
$$z_i = \frac{1}{k} \sum_{j=1}^k y_{i+j-1}, \quad i = \overline{1, (m-1)k+1}$$

## Obs:

- Numărul de elemente din serie este redus de la  $n$  la  $n-k+1$
- Variațiile pe termen scurt pot fi pierdute prin mediere
- Mediarea poate fi unidirecțională (se utilizează doar valori anterioare momentului curent) sau bidirecțională/ centrată (se utilizează atât valori anterioare cât și ulterioare)

# Pre-procesarea seriilor de timp

Exemplu (Moving average smoothing)



Obs:

- cu cât este mai mare dimensiunea ferestrei ( $k$ ) cu atât este mai semnificativă netezirea (se pot pierde inclusiv variații relevante)
- la calculul mediei pot participa toate elementele (inclusiv cele îndepărtate în timp) dar cu influență mai mică → **netezire exponențială**

# Pre-procesarea seriilor de timp

## Netezire exponențială

**Idee:** valoarea netezită se definește ca o combinație liniară a valorii curente și a valorii netezite anterioare

$$z_i = \alpha \cdot y_i + (1 - \alpha) \cdot z_{i-1}, \quad i = \overline{1, n}$$

$$z_0 = y_1$$

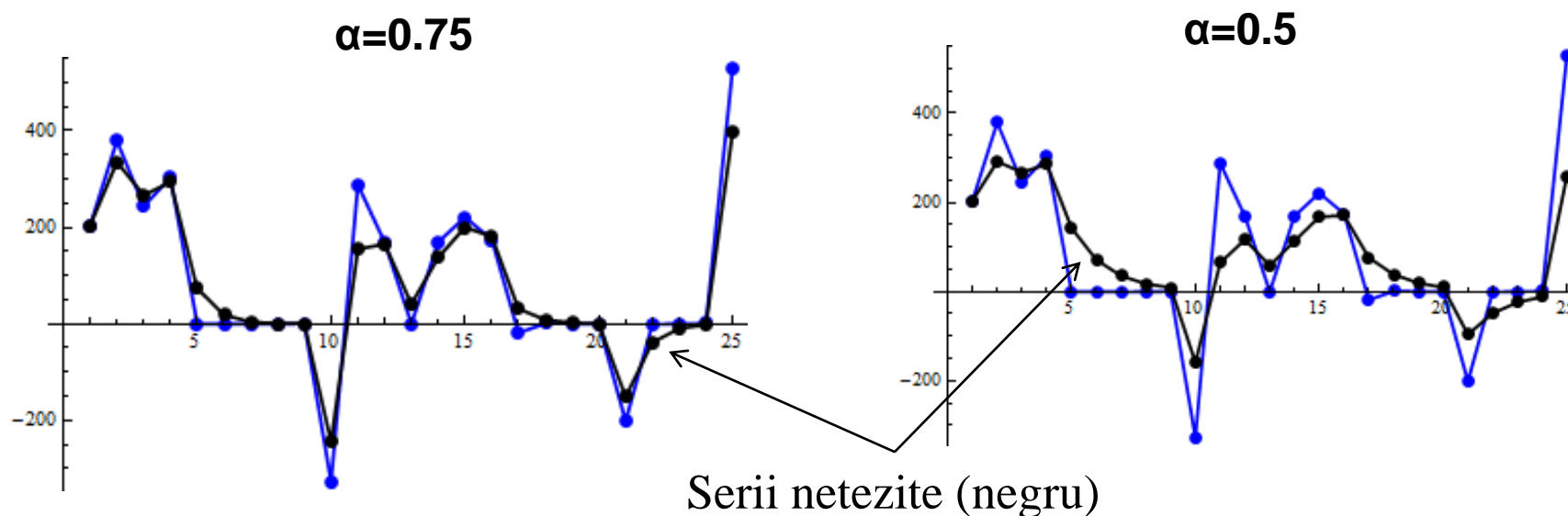
$$z_i = (1 - \alpha)^i z_0 + \alpha \sum_{j=1}^i y_j (1 - \alpha)^{i-j}, \quad i = \overline{1, n}$$

## Obs:

- Dacă  $\alpha=1$  atunci nu se aplică netezire; dacă  $\alpha=0$  toată seria este netezită (va avea valoarea primului element)
- Netezirea exponențială se bazează pe ideea că valorile mai recente sunt mai importante iar cele mai “vechi” au o influență mai mică; **influența valorilor anterioare este controlată prin  $1-\alpha$**

# Pre-procesarea seriilor de timp

Exemplu (exponential smoothing)



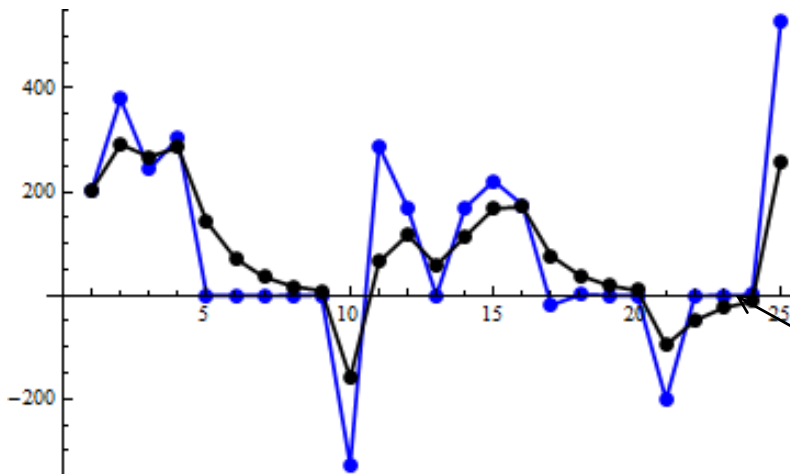
Obs:

- cu cât  $\alpha$  este mai mare (mai aproape de 1) cu atât seria este mai puțin netezită
- cu cât  $\alpha$  este mai mică (mai aproape de 0) cu atât seria este mai semnificativ netezită

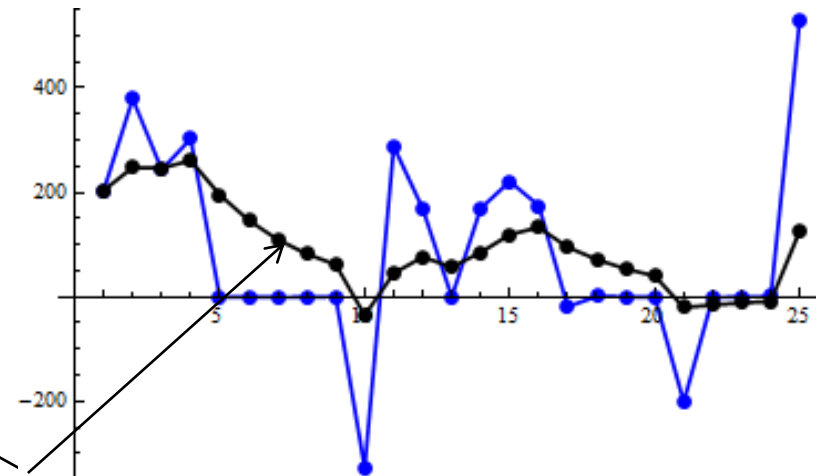
# Pre-procesarea seriilor de timp

Exemplu (exponential smoothing)

$\alpha=0.5$



$\alpha=0.25$



Serii netezite (negru)

Obs:

- cu cât  $\alpha$  este mai mare (mai aproape de 1) cu atât seria este mai puțin netezită
- cu cât  $\alpha$  este mai mică (mai aproape de 0) cu atât seria este mai semnificativ netezită

# Pre-procesarea seriilor de timp

## Scalare și standardizare

### Variante:

#### Scalare:

$$z_i = \frac{y_i - \min(y)}{\max(y) - \min(y)}, \quad i = \overline{1, n}$$

#### Standardizare:

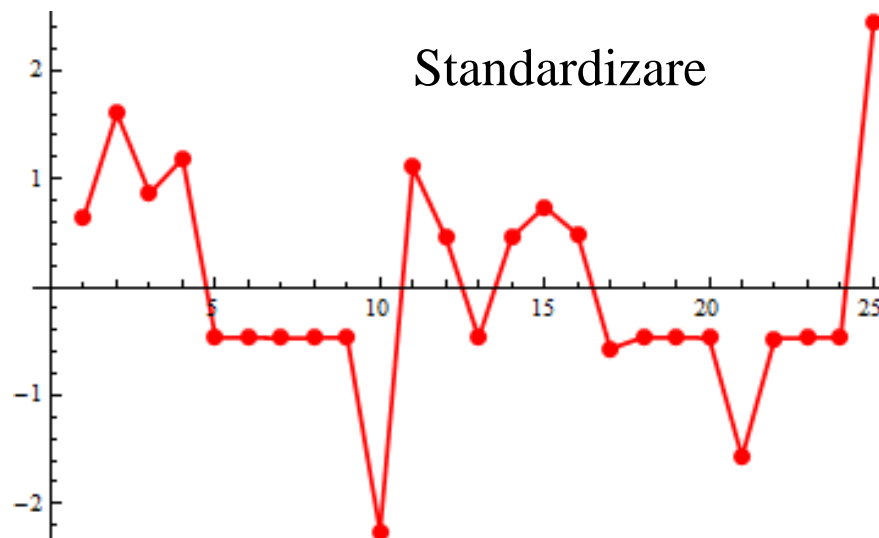
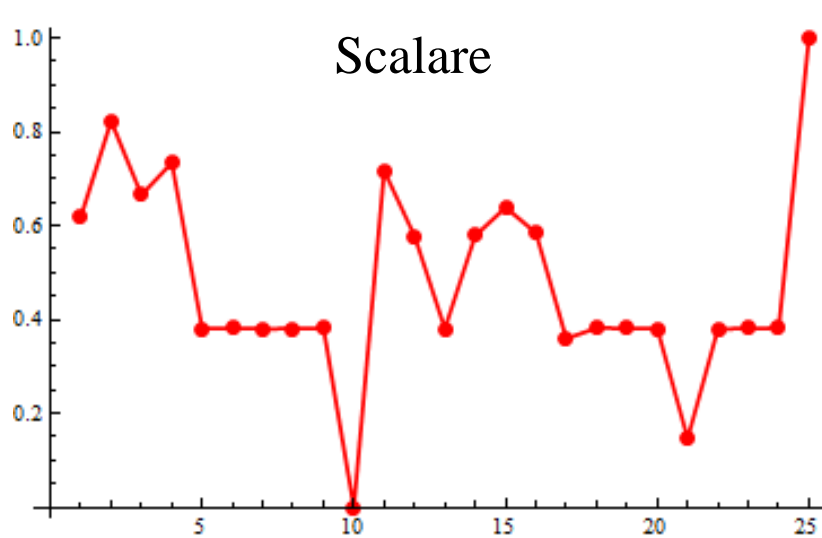
$$z_i = \frac{y_i - \text{mean}(y)}{\text{stdev}(y)}, \quad i = \overline{1, n}$$

#### Obs:

- $\min(y)$  și  $\max(y)$  reprezintă valoarea minimă respectiv cea maximă din serie
- $\text{mean}(y)$  și  $\text{stdev}(y)$  sunt valoarea medie respectiv abaterea standard

# Pre-procesarea seriilor de timp

## Exemplu



## Obs:

- Scalarea și standardizarea conservă forma seriei dar schimbă domeniul de valori

# Predicție

## Scop:

- Estimarea prețului viitor al unei acțiuni, prognoza meteo, estimarea evoluției unor indicatori economici etc

## Predicție:

- **Intrare:** una sau mai multe serii de timp
- **Ieșire:** valori viitoare ale seriei

## Cum poate fi abordată problema?

- Varianta 1: ca o problemă de **regresie** – se estimează explicit dependența dintre attributele comportamentale și timp (necesită cunoașterea explicită a valorilor variabilei contextuale = momentele de timp) sau se exploatează relația dintre două serii de timp
- Varianta 2: **utilizând modele care exprimă relația dintre valori curente și valori anterioare ale seriei (modele autoregresive)** – corelația dintre valori ale seriei este denumită de regula **autocorelație**



# Analiza seriilor de timp

## Serii staționare

- Intuitiv, o serie staționară se caracterizează prin faptul că proprietățile sale statistice (medie, varianță, autocorelație = corelație între valori din serie) sunt constante în timp

## Staționaritate strictă:

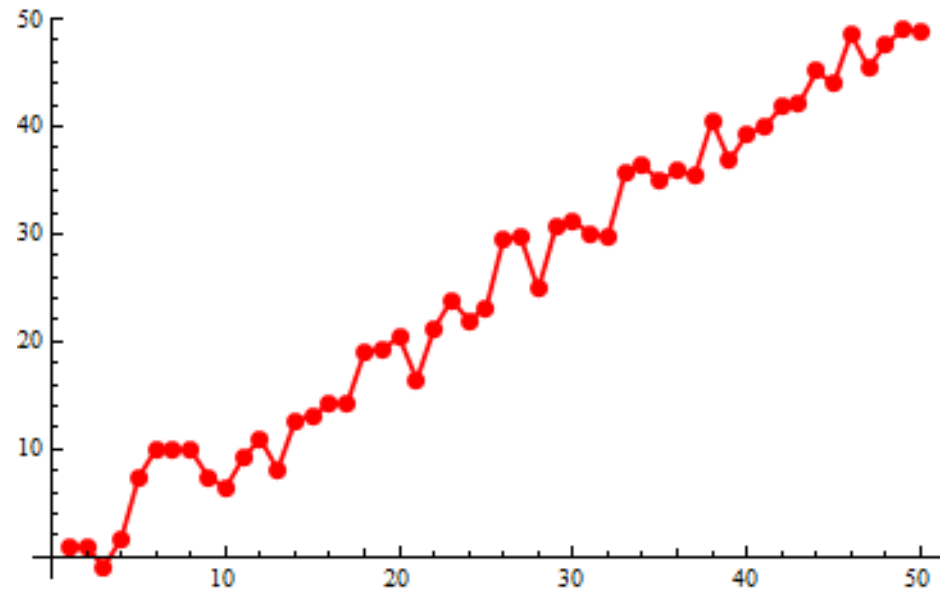
- distribuția de probabilitate a valorilor din orice interval de timp  $[a,b]$  este identică cu distribuția de probabilitate a valorilor din intervalul translatat  $[a+h, b+h]$  (pentru un  $h>0$  arbitrar)

## Obs:

- Cum se verifică staționaritatea? Se estimează indicatori statistici bazat pe ferestre de timp și se analizează dacă se obțin valori similare pt ferestre diferite
- În cazul seriilor nestaționare înainte de a aplica o tehnică de predicție autoregresivă ar fi util ca seria să fie transformată într-una staționară sau să se folosească o tehnică care țină cont de faptul că seria e nestaționară.

# Analiza seriilor de timp

Exemplu: serie artificial construită:  $y_i = i + \text{zgomot}$  (zgomotul este generat folosind o distribuție normală de medie 0 și abatere standard 2)



# Analiza seriilor de timp

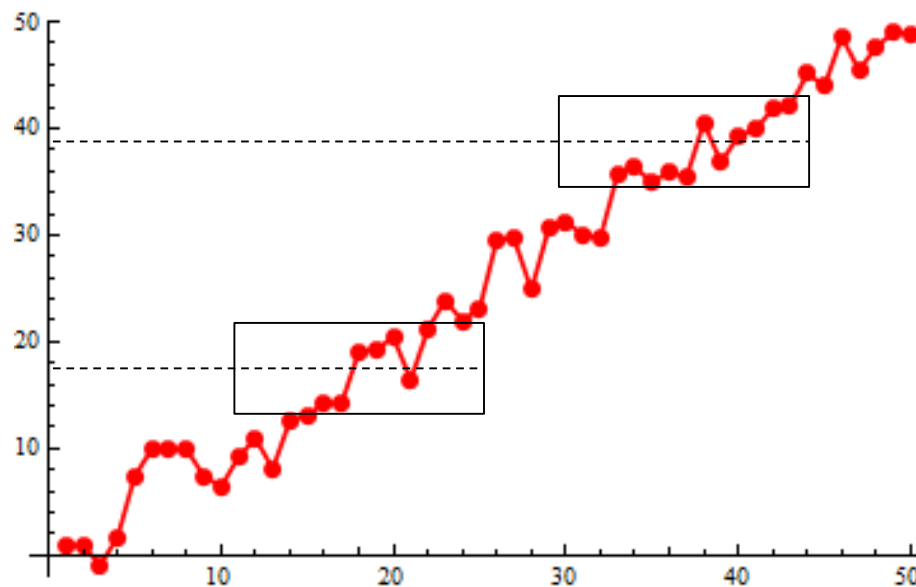
**Exemplu:** serie artificial construită:  $y_i = i + \text{zgomot}$  (zgomotul este generat folosind o distribuție normală de medie 0 și abatere standard 2)

**Obs:**

- Aceasta este o **serie nestaționară** întrucât mediile valorilor corespunzând unor ferestre de timp diferite sunt diferite

Medie 2  
(a doua fereastră)

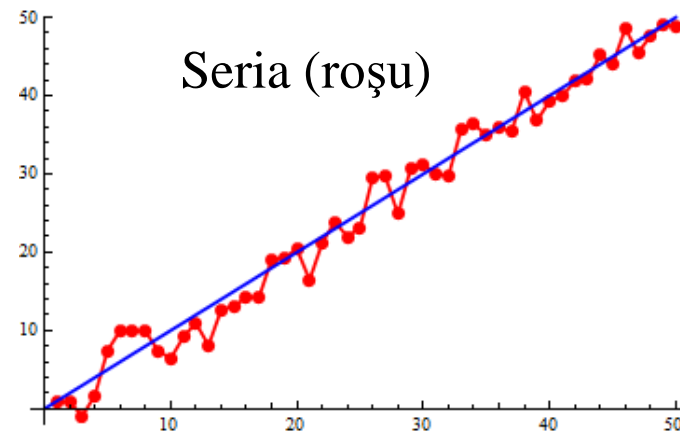
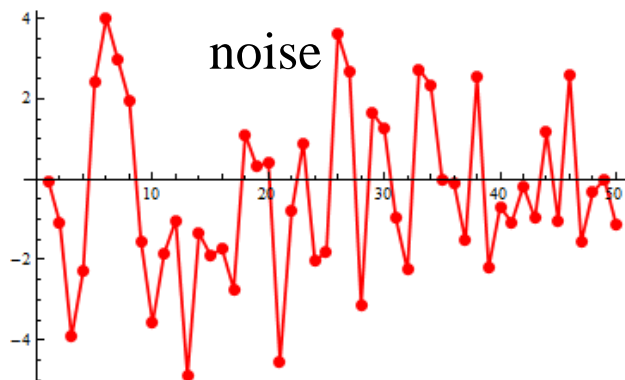
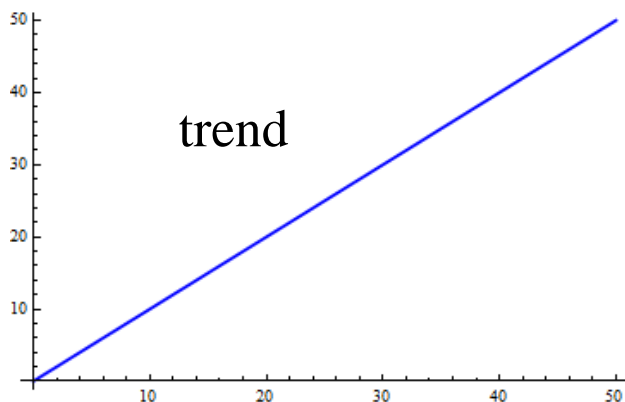
Medie 1  
(prima fereastră)



# Analiza seriilor de timp

**Exemplu:** serie artificial construită:  $y_i = i + \text{zgomot}$  (zgomotul este generat folosind o distribuție normală de medie 0 și abatere standard 2)

**Obs.** Sunt 2 componente: **tendința (trend)** și **zgomot (noise)**

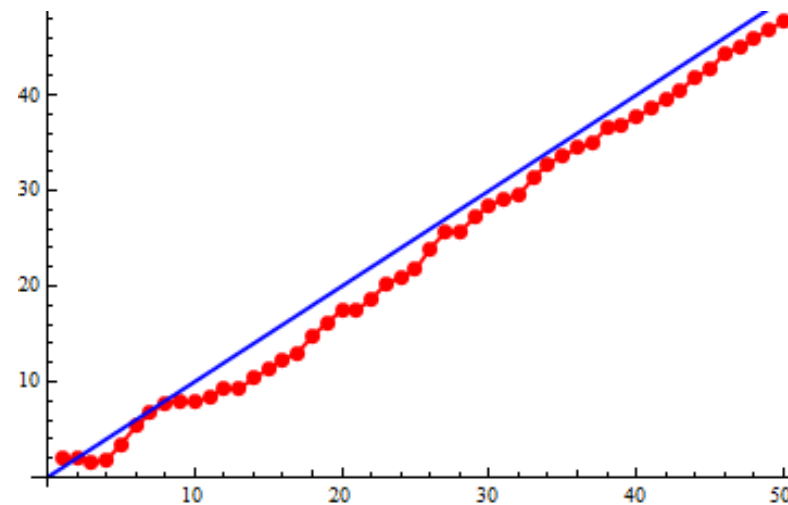
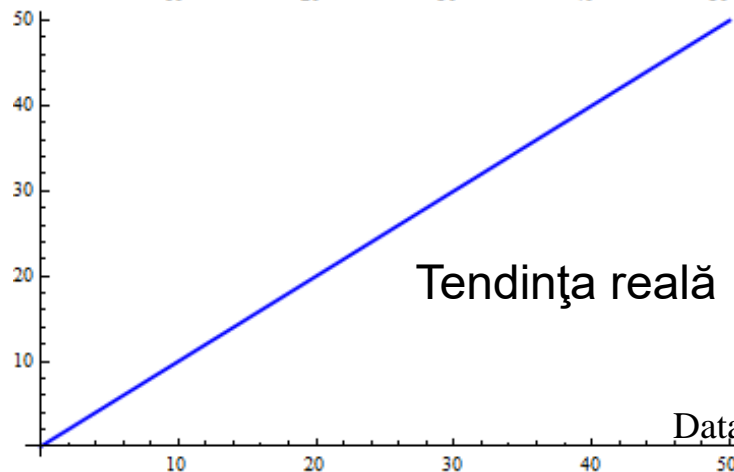
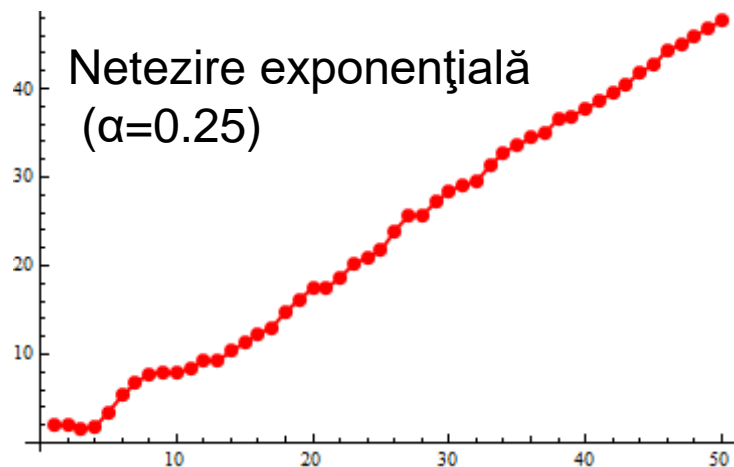


Cum pot fie extrase cele două componente din seria inițială?

# Analiza seriilor de timp

Extragerea tendinței = eliminarea zgomotului

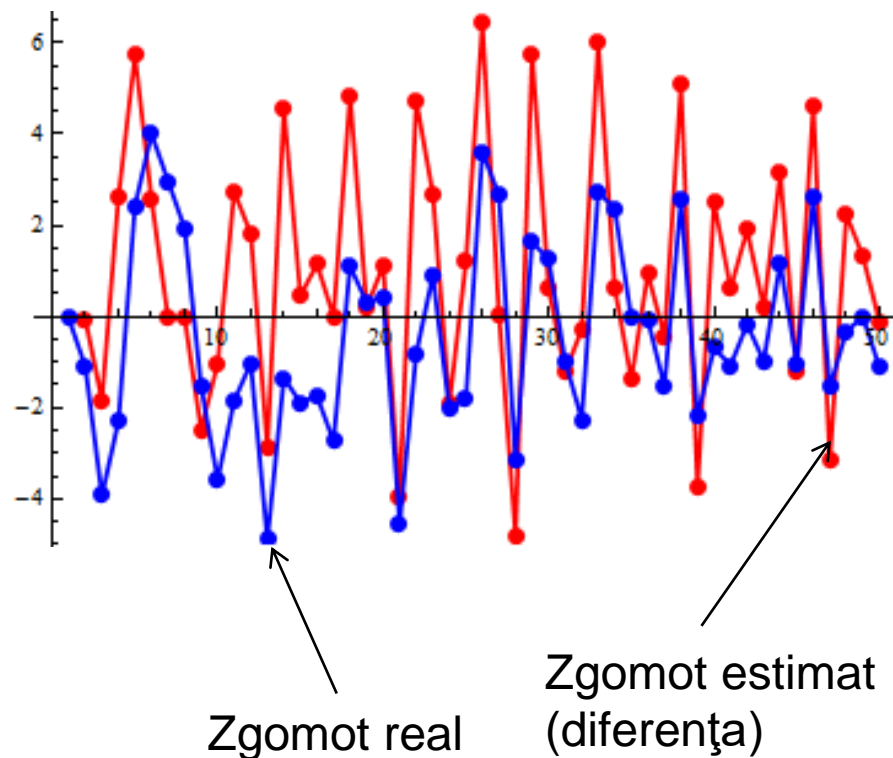
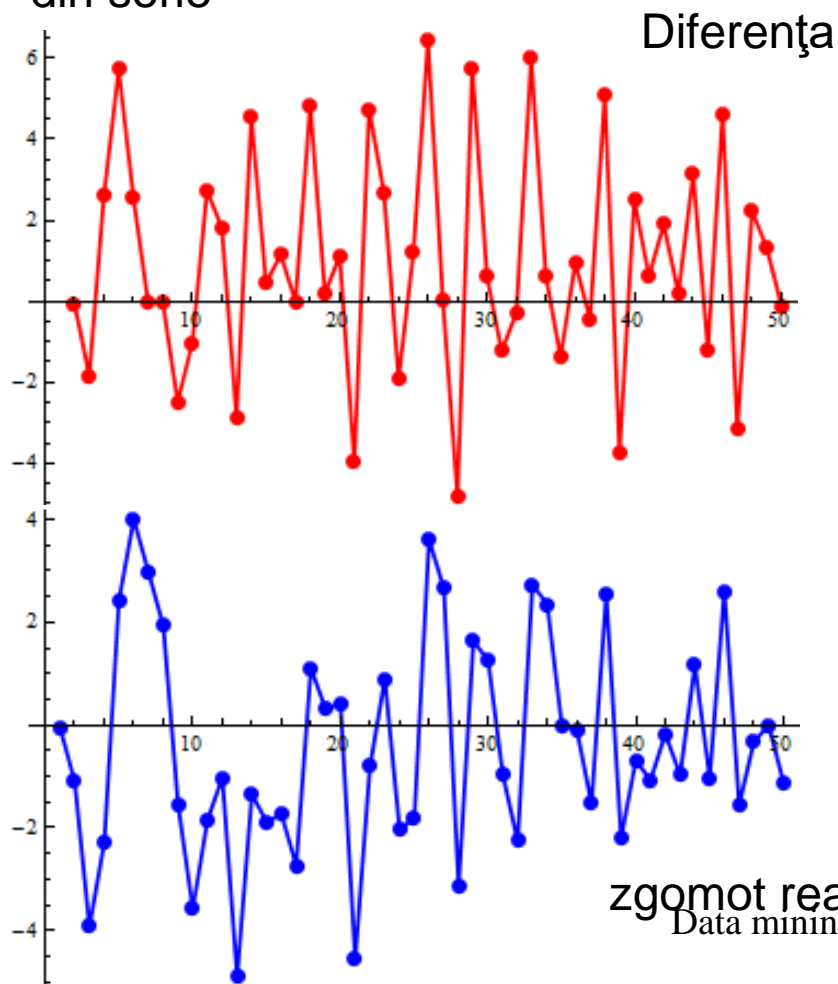
Cum poate fi realizată? prin netezire



# Analiza seriilor de timp

Extragerea zgomotului (fluctuației) = eliminarea tendinței

Cum poate fi realizată? prin calculul diferenței dintre elementele succesive din serie



# Analiza seriilor de timp

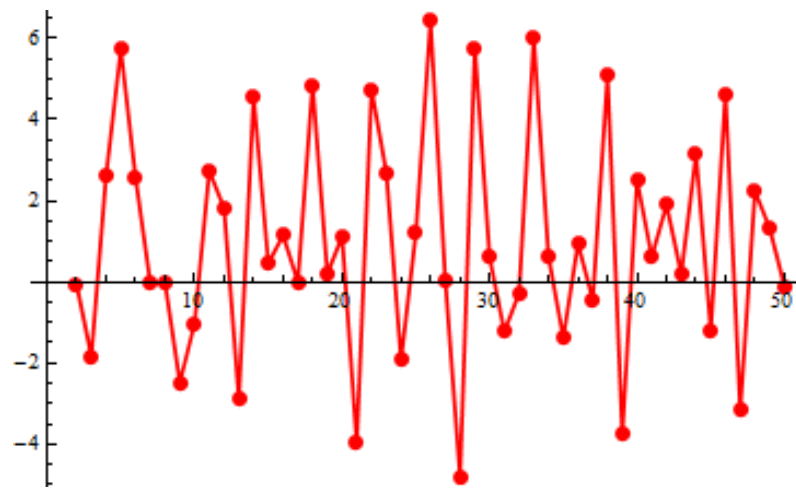
Extragerea zgomotului (fluctuației) = eliminarea tendinței

Cum poate fi realizată? prin calculul diferenței dintre elementele succesive din serie

Transformare prin calculul diferenței:  $z_i = y_i - y_{i-1}$

Obs:

- Seria obținută prin diferențiere este staționară



Diferența

# Analiza seriilor de timp

Extragerea zgomotului (fluctuației) = eliminarea tendinței

Alte variante:

- Eliminarea efectului sezonier

$$z_i = y_i - y_{i-T}$$

- La seriile cu creștere geometrică (de exemplu serii de prețuri în care factorul de inflație e constant) poate fi utilă **logaritmare** înainte de calculul diferențelor

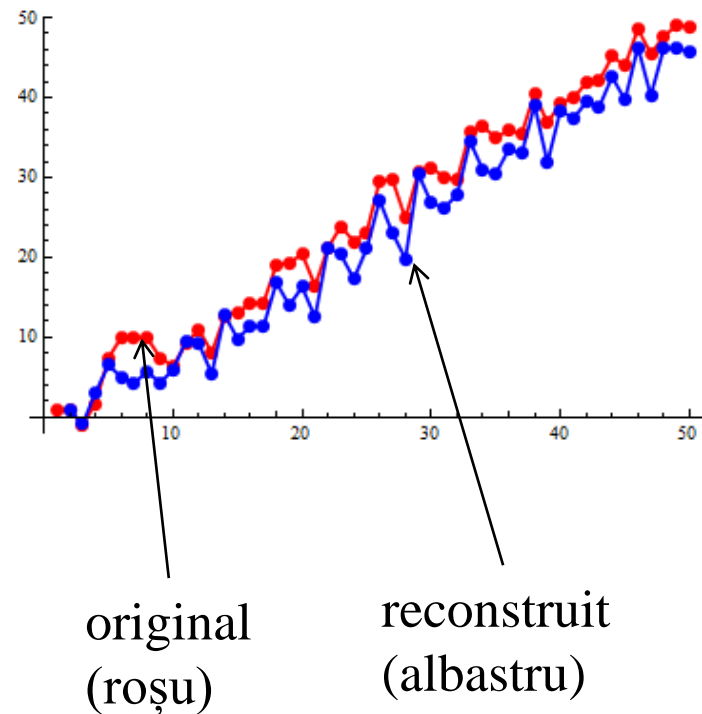
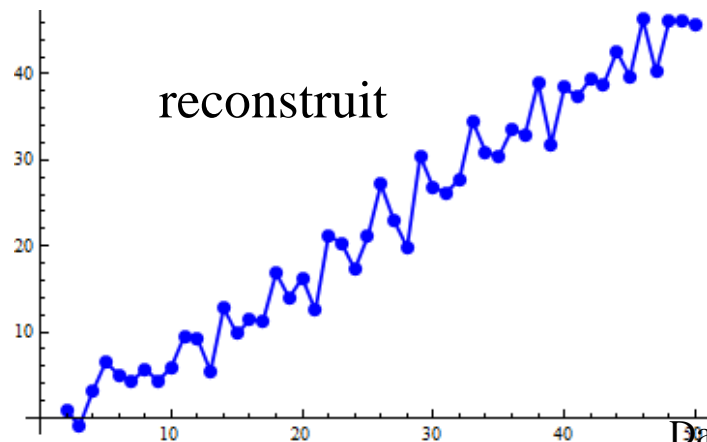
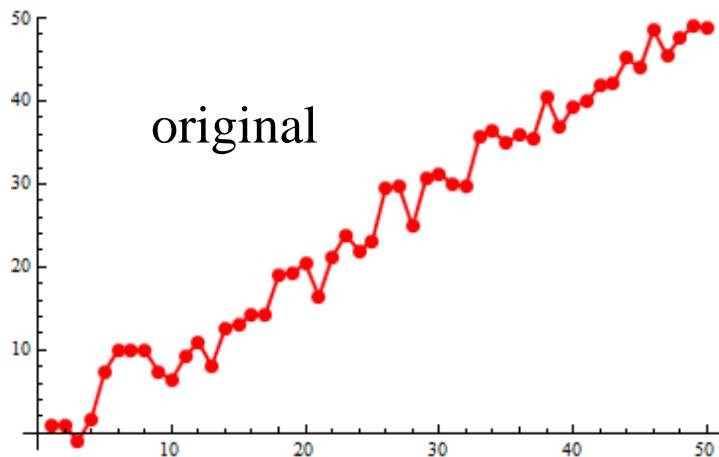
Întrebare:

- Poate fi reconstruită seria inițială pornind de la estimările tendinței și zgomotului?



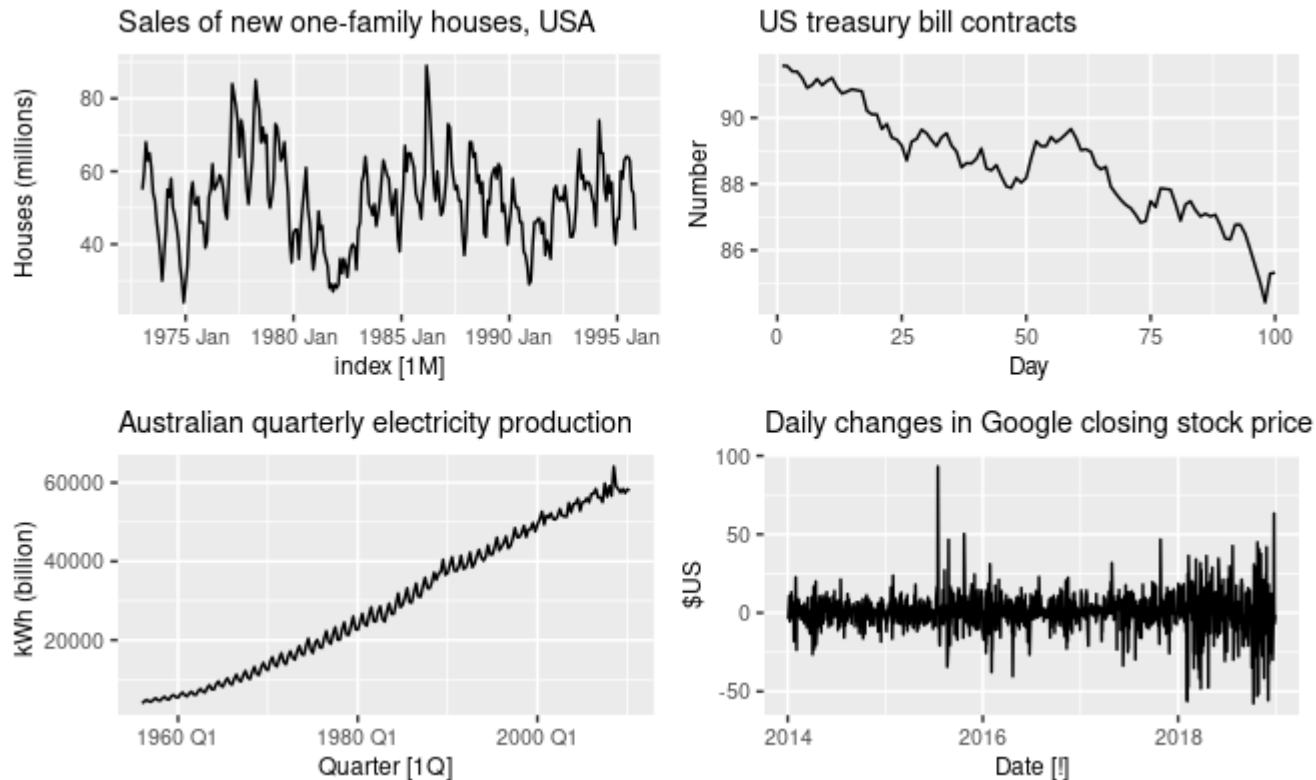
# Tendință (trend) și zgomot (noise)

Reconstruire: suma dintre estimarea tendinței și estimarea zgomotului



# Tendință, caracter sezonier, zgomot

Exemple de serii cu diferite comportamente:



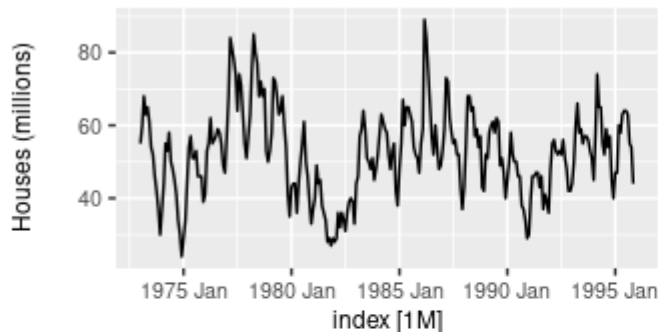
Sursa: <https://otexts.com/fpp3/tspatterns.html>

# Tendință, caracter sezonier, zgomot

Exemple de serii cu diferite comportamente:

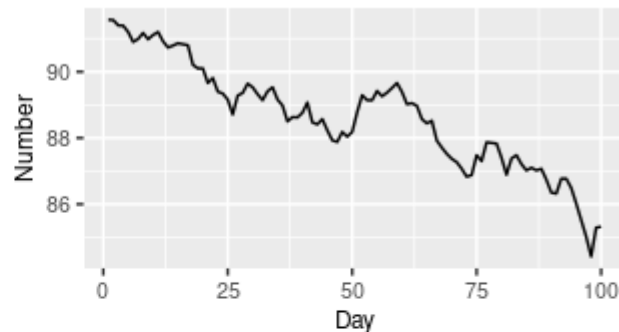
1

Sales of new one-family houses, USA



2

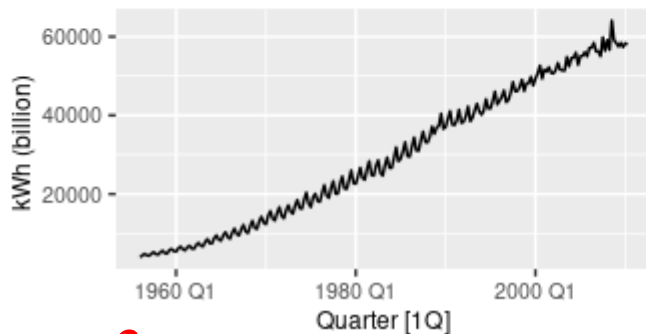
US treasury bill contracts



1. Fără tendință dar cu caracter sezonier

2. Tendință descrescătoare, fără caracter sezonier

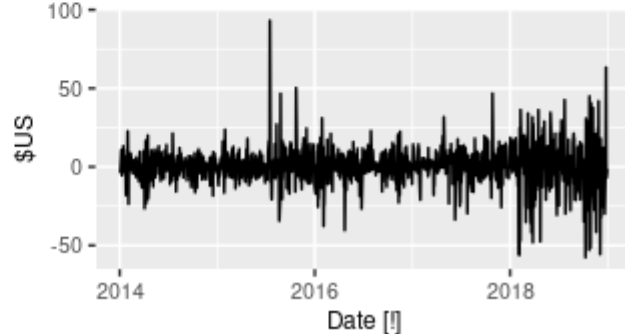
Australian quarterly electricity production



3

4

Daily changes in Google closing stock price



3. Tendință crescătoare cu caracter sezonier prezent

4. Fără tendință sau caracter sezonier prezent (doar fluctuații aleatoare)

# Predicție

Intrare:  $y(1), y(2), \dots, y(t)$

Ieșire:  $\hat{y}(t + h), h \geq 1$

Variante:

- **Regresie** (în raport cu variabila contextuală sau în raport cu valorile unei alte serii):  
 $\hat{y}(t) = F(t)$  sau  $\hat{y}(t) = F(x(t))$
- **Predicție naivă:**
  - $\hat{y}(t + h) = y(t)$  (se ia în considerare ultima valoare din serie)
  - $\hat{y}(t + h) = y(t + h - T)$  (se ține cont de caracterul sezonier – perioada = T)
- **Metoda mediei:**  $\hat{y}(t + h) = (y(t) + y(t - 1) + \dots + y(t - k + 1))/k$
- **Modele bazate pe descompunere și netezire (exponential smoothing):**
  - serie = trend + componenta sezoniera + zgomot (model aditiv)
  - serie = trend \* componenta sezoniera \* zgomot (model multiplicativ – poate fi transformat în model aditiv prin logaritmare)
- **Modele autoregresive:**  $\hat{y}(t + 1) = F(y(t), y(t - 1), \dots, y(t - p + 1))$

# Metoda netezirii

Cum poate fi estimată o nouă valoare din serie?

- Idee: utilizarea netezirii exponențiale în locul mediei simple:

$$\hat{y}(t + 1) = \alpha y(t) + (1 - \alpha) \hat{y}(t), \alpha \in [0, 1]$$

- **Dezavantaj:** nu include prezența unui trend

**Soluție:** metoda Holt Winters

- se bazează pe ideea estimării iterative a nivelului de bază ( $a(t)$ ) și a unei pante ( $b(t)$ ):

$$\hat{y}(t + h) = a(t) + hb(t)$$

$$a(t) = \alpha y(t) + (1 - \alpha)(a(t - 1) + b(t - 1))$$

$$b(t) = \beta(a(t) - a(t - 1)) + (1 - \beta)b(t - 1)$$

# Metoda netezirii

Extindere: includere componenta sezonieră

- Se bazează pe ideea estimării iterative a componentelor folosind netezirea exponențială (nivel mediu, panta tendinței, componenta sezonieră)

$$y(t + h) = a(t) + h \cdot b(t) + s(t - T + 1 + (h - 1) \bmod T)$$

$$a(t) = \alpha \cdot (y(t) - s(t - T)) + (1 - \alpha) \cdot (a(t - 1) + b(t - 1))$$

$$b(t) = \beta \cdot (a(t) - a(t - 1)) + (1 - \beta) \cdot b(t - 1)$$

$$s(t) = \gamma \cdot (y(t) - a(t)) + (1 - \gamma) \cdot s(t - T)$$

Obs:

- parametrii  $\alpha, \beta, \gamma$  se determină prin minimizarea erorii medii pătratice pe valorile din serie (toate valorile sunt în  $(0,1)$ )
- În cazul unei serii staționare este suficient să se estimeze, prin netezire exponențială, componenta  $a(t)$  corespunzătoare nivelului mediu

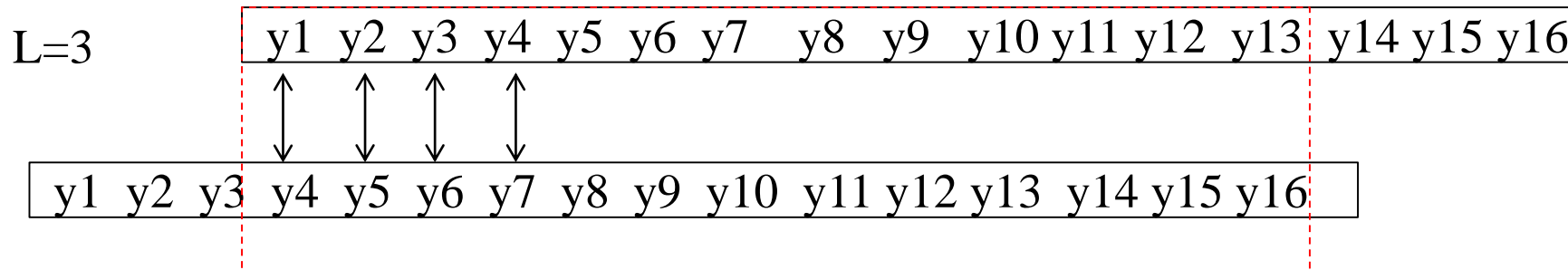
# Modele autoregresive

Ideea de bază:

- dacă valoarea **autocorelației** (corelație între valorile aflate în serie) este mare (în valoare absolută) atunci valoarea corespunzătoare unui moment poate fi estimată pe baza valorilor din vecinătate

Autocorelație pt o serie staționară,  $(y_1, y_2, \dots, y_n)$  = corelația dintre valori separate prin întârzierea  $L$  (lag)

$$Autocorrelation(y_i, y_{i+L}) = \frac{\frac{1}{n-L} \sum_{i=1}^{n-L} (y_i - \text{avg}(Y))(y_{i+L} - \text{avg}(Y))}{\text{var}(Y)}$$



# Modele autoregresive

**Autocorelație (acf = autocorrelation function):** măsură a corelației între valorile corespunzătoare unor momente de timp aflate la aceeași distanță

$$acf(y_i, y_{i+L}) = \frac{\frac{1}{n-L} \sum_{i=1}^{n-L} (y_i - \text{avg}(Y))(y_{i+L} - \text{avg}(Y))}{\text{var}(Y)}$$

**Autocorelație parțială:** măsură a corelației între valorile corespunzătoare unor momente de timp aflate la aceeași distanță din care **s-a eliminat dependența indusă de valorile intermediare**

$$pacf(y_i, y_{i+L}) = acf(y_i - z_i, y_{i+L} - z_{i+L})$$

unde

- $z_i$  este o estimare a lui  $y_i$  pe baza unui model de regresie liniară în raport cu  $y_{i+1}, y_{i+2}, \dots, y_{i+L-1}$  ( $z_i = b_1 y_{i+1} + b_2 y_{i+2} + \dots + b_{L-1} y_{i+L-1}$ )
- $z_{i+L}$  este o estimare a lui  $y_{i+L}$  pe baza unui model de regresie liniară în raport cu  $y_{i+L-1}, y_{i+L-2}, \dots, y_{i+1}$  ( $z_{i+L} = c_1 y_{i+L-1} + c_2 y_{i+L-2} + \dots + c_{L-1} y_{i+1}$ )



# Modele autoregresive

Forma generală a unui model autoregresiv de ordin  $p$ : AR( $p$ )

$$y_t = \sum_{i=1}^p a_i y_{t-i} + c + \varepsilon_t$$

Obs:

- $p$  este ordinul modelului și poate fi ales analizând diferite valori posibile ale întârzierii  $L$ :
  - $p$  se alege ca fiind prima valoare  $L$  (pornind de la  $L=1$  și crescând valoarea lui  $L$ ) pt care valoarea absolută a auto-corelației parțiale (pacf) este suficient de mică
- $a_1, a_2, \dots, a_p$  și  $c$  sunt parametri ai modelului și se estimează folosind date de antrenare și metoda celor mai mici pătrate (similar cu modelele de regresie liniară)
- $\varepsilon_t$  reprezintă componenta de tip zgomot (care nu e modelată explicit în modelele de tip AR)

# Modele autoregresive

Modele de tip medie mobilă (Moving Average): MA(q)

Motivație:

- Modelele autoregresive simple nu pot explica toate variațiile (în mod particular schimbările bruște, de tipul șocurilor)

Idee:

- Modelele de tip MA prezic valorile următoare pe baza **deviațiilor anterioare ale valorilor reale față de cele prezise** (valoarea din serie este determinată de un eveniment aleator care poate avea impact nu doar în momentul curent ci și în câteva momente următoare)

$$y_t = \sum_{i=1}^q b_i \varepsilon_{t-i} + c + \varepsilon_t,$$

cu  $\varepsilon_t$  variabile aleatoare independente și identic distribuite

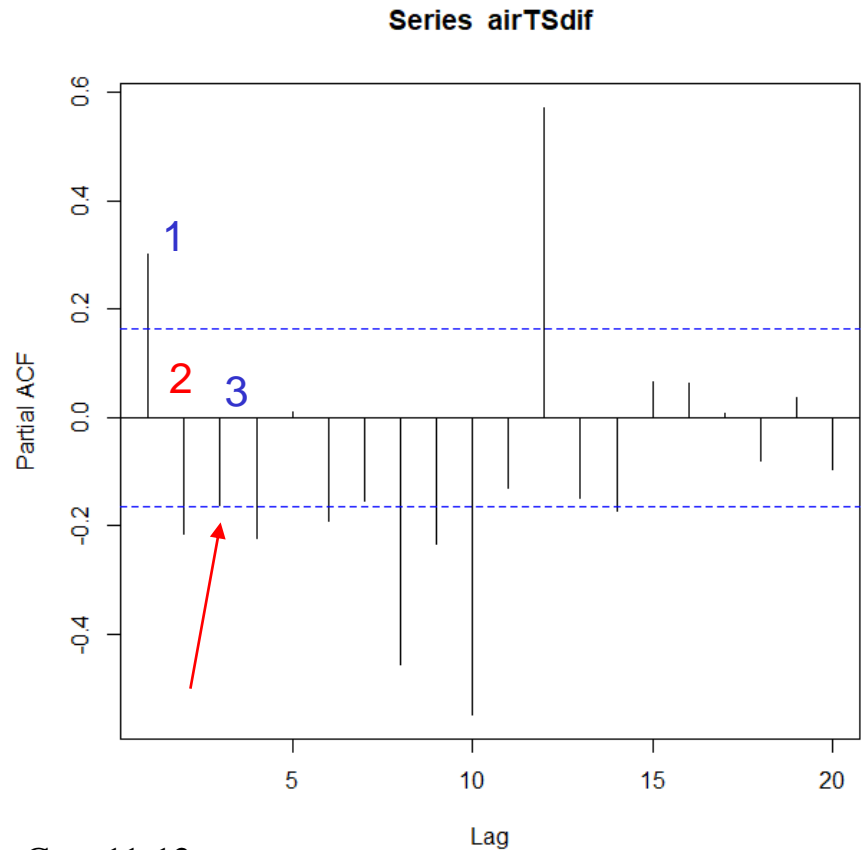
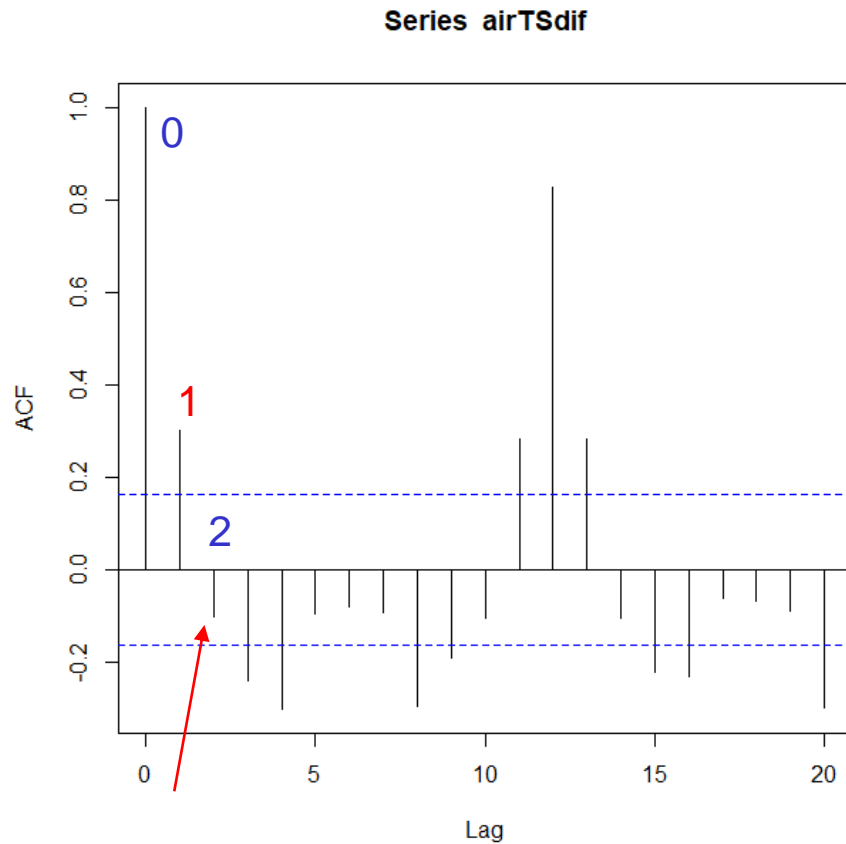
Obs:

- Presupunând că **seria este staționară și zgomotul are medie 0**, valoarea lui **c este de fapt media valorilor din serie**
- Parametrii  $b_1, b_2, \dots, b_q$  se estimează din date
- Parametrul  $q$  se alege ca fiind numărul de valori ale întârzierii până la prima cu o valoare absolută a **auto-corelației (acf)** suficient de mică

# Modele autoregresive

Autocorelație (acf):  $q=1$

Autocorelație parțială (pacf):  $p=2$



# Modele autoregresive

Modele autoregresive combinate: ARMA(p,q)

Motivație:

- Se combină capacitatea de predicție a modelelor autoregresive și a celor bazate pe medie mobilă:

$$y_t = \sum_{i=1}^p a_i y_{t-i} + \sum_{i=1}^q b_i \varepsilon_{t-i} + c + \varepsilon_t$$

Obs:

- Un aspect important este alegerea valorilor p și q: se determină analizând valorile autocorelației și ale autocorelației parțiale
- În cazul în care seria nu este staționară se poate folosi modelul **ARIMA(p,d,q)** (**Autoregressive Integrated Moving Average**) care asigură eliminarea tendinței prin aplicarea unui operator “diferența” de ordin **d** (d se stabilește ca fiind numărul de iterații ale operatorului “diferența” până când se obține o serie staționară)

# Descoperirea șabloanelor

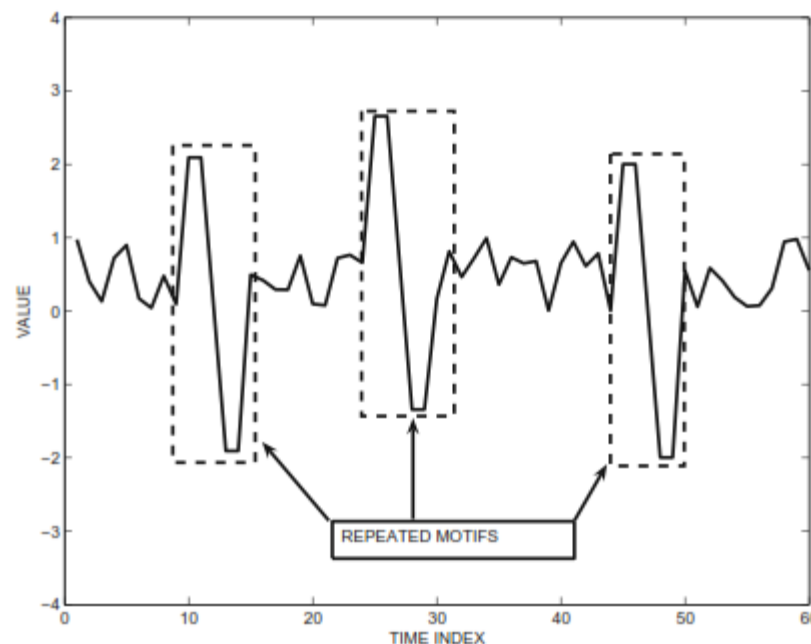
Șablon (motiv) = structură ce apare frecvent în serie

Procesul de descoperire

Intrare:

- Cel puțin o serie
- Lungimea  $L$  a șablonului
- Măsură de similaritate/ disimilaritate
- Prag pentru similaritate/ disimilaritate

Ieșire: subsecvență de lungime  $L$  ce apare frecvent în serie



C. Aggarwal, Data Mining – the Textbook, 2015

# Descoperirea șabloanelor

Șablon (motiv) = structură ce apare frecvent în serie

Exemplu: Algoritm de tip forță brută

FindMotif ( $y[1..n]$ ,  $L$ ,  $\epsilon$ )

countMax=0

FOR  $i=1, n-L+1$  DO

    candidate= $y[i..i+L-1]$

    count=0

    FOR  $j=1, n-L+1$  DO

$D = \text{dist}(y[i..i+L-1], y[j..j+L-1])$

        IF ( $i \neq j$ ) and ( $D \leq \epsilon$ ) THEN count=count+1

    ENDFOR

    IF count[i]>countMax THEN best=i; countMax=count

ENDFOR

RETURN ( $y[\text{best}..\text{best}+L-1]$ )

# Excepții (anomalii)

Există două tipuri de excepții (anomalii) într-o serie de date:

## Excepții (anomalii) punctuale:

- Deviație semnificativă de la valoarea prezisă
- Corespunde unei schimbări bruște în seria de date

## Excepții (anomalii) în privința formei:

- O succesiune de valori poate reprezenta o anomalie chiar dacă valorile individuale nu sunt neobișnuite
- De exemplu, într-o electrocardiogramă o bătaie neregulată a inimii poate fi considerată o anomalie

# Excepții (anomalii)

## Detecția anomaliilor punctuale:

Step 1: se determină valoarea prezisă (pe baza modelului construit folosind valorile anterioare)  $(z_m, z_{m+1}, \dots, z_n)$

Step 2: se construiește seria deviațiilor  $(d_m, d_{m+1}, \dots, d_n)$  cu  $d_i = z_i - y_i$

Step 3: se calculează deviațiile standardizate  $(s_m, s_{m+1}, \dots, s_n)$  cu  $s_i = (d_i - \text{avg}(d)) / \text{stdev}(d)$

Dacă valoarea absolută a lui  $s_i$  este mai mare decât un prag (e.g. 3) atunci se consideră că este anomalie



# Excepții (anomalii)

## Detecția anomaliilor de formă:

Step 1: se extrag toate subseriile corespunzătoare unui ferestre de dimensiune  $W$

Step 2: se calculează **distanța (disimilaritatea)** dintre fiecare subserie și toate celelalte corespunzătoare unor ferestre disjuncte

Step 3: Subseriile care diferă semnificativ de celelalte sunt considerate excepții potențiale

## Probleme:

- Alegerea lui  $W$
- Alegerea pragului

# Măsurarea (di)similarității

- Serii “aliniate” (valorile din cele două serii corespund acelorasi momente de timp)
  - Se poate utiliza orice măsură de (di)similaritate corespunzătoare unor date vectoriale
- Seriile nu sunt aliniate (de exemplu sunt înregistrări audio cu viteze diferite)
  - Se poate folosi un algoritm de matching între serii – similar algoritmilor de aliniere de secvențe biologice
  - **Dynamic Time Warping** (bazat pe tehnica programării dinamice)

## Idee DTW:

- $D(i,0)=D(0,j)=\inf, \quad i=1..n, \quad j=1..m$
- $D(0,0)=0$
- For  $i=1..n$ 
  - For  $j=1..m$ 
    - $D(i,j)=\text{dif}(x[i],y[j])+\min\{D(i-1,j), D(i,j-1), D(i-1,j-1)\}$
- $\text{Dist}(x[1..n],y[1..m]) = D(n,m)$

# Sumar

- Pre-procesarea seriilor temporale:
  - Imputarea valorilor absente folosind interpolare
  - Normalizare/standardizare
  - Netezire (pentru eliminarea zgomotului); ex: exponential smoothing
- Descompunerea unei serii în componente:
  - Tendința – obținută prin netezire
  - Zgomot – obținut prin eliminare tendința - calcul diferențe
    - ordin 1 pentru eliminare tendință liniară
    - ordin 2 pentru tendință pătratică
- Predicție
  - Bazată pe netezire exponențială
  - Modele autoregresive si de tip moving average (AR, MA)
  - Modele combinate (ARMA, ARIMA)
- Biblio adițional: <https://otexts.com/fpp3/>

# Cursuri următoare

- Metode de tip ansamblu (meta-modele)
- Colecții de modele (bucket of models)
- Colecții de arbori aleatori (random forests)
- Strategii de agregare a modelelor
  - Bagging
  - Boosting
  - Stacking