

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

1 Curs 5. Tipuri de date. (descriptive

In [17]:

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

print(f'NumPy version: {np.__version__}')
print(f'pandas version: {pd.__version__}')
```

executed in 5.42s, finished 19:45:20 2021-03-21

NumPy version: 1.19.2
pandas version: 1.2.3

1.1 Tipuri de date

- Un set de date este o colectie de obiecte-data
- Sinonime pentru obiecte-data: înregistrare, pu confuzie), observatie, entitate.
- Obiectele sunt descrise prin attribute
- Sinonime pentru atribut: **variabila**, **caracteris**

Definitie: Atribut: proprietate sau caracteristica a u

- Exemple: culoarea ochilor, temperatura.
- Trebuie facuta diferenta între proprietatile atrib
 - acelasi atribut poate avea valori diferite: în
 - diferite attribute pot fi masurate cu acelasi reprezentate prin numere întregi; în timp

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transform
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Seturi
 - 1.1.4.2 Seturi
 - 1.1.4.3 Seturi
 - ▼ 1.2 Calitatea datelor
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transform
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

1.1.1 Tipuri de attribute

Există diferite tipuri de attribute:

- Categoriale (calitative)
 - **nominale**: valori diferite care permit recu
 - **ordinale**: valorile permit ordonarea obiec
- operatii permise: =, \neq , <, >; functii aplicat
- Numerice (cantitative)
 - **interval**: se poate face diferenta între val
 - Fahrenheit; pe lângă operatiile de mai sus
 - **multiplicabile**: permit împărțiri și înmultir
 - mai sus și înmultirea, împărțirea; functii a

1.1.2 Transformari la nivel de atribut

Există niste transformari care se pot efectua sau n

- pentru attribute nominale: orice asociere unu-l
- datelor;
- pentru attribute ordinale: orice modificare de v
- monoton crescatoare; {bun, mai bun, cel mai l
- pentru attribute interval: transformari de forma
- pentru attribute multiplicabile: val_veche/val_n

1.1.3 Descrierea atributelor prin num

- Attribute discrete:
 - o multime cel mult numarabila de valori;
 - exemple: coduri postale, cuvinte într-un d
 - se reprezinta cel mai frecvent ca numere
 - caz special — attribute binare: {prezent, a
- Attribute continue:
 - valorile sunt exprimate prin numere reale
 - exemple: temperatura, masa
 - dpdv practic reprezentarea se face cu o p
 - reprezentare actuala: valori în virgula mol
- Valori asimetrice:
 - doar prezenta unei trasaturi (i.e. valoare r
 - exemple: vectorul care reprezinta daca ni
 - daca se iau in considerare doi astfel de v
 - simultan

1.1.4 Tipuri de seturi de date

- Seturi de date de tip: înregistrare, de tip grafic
- Caracteristici generale:
 - **dimensionalitatea** = numărul de atribute
multe dimensiuni se pot aplica tehnici de
 - **caracterul rarefiat al datelor** = procentul
poate reduce drastic necesarul de memorie
 - **rezoluția** = scara la care se face raportarea
zile pot arăta iminența unei furtuni, dar la

1.1.4.1 Seturi de date de tip înregistrare

- cel mai des furnizate și frecvent utilizate în aplicații
- nu există legătură între înregistrări distincte
- stocare: fișiere text (e.g. CSV), Excel, baze de date

Cazuri remarcabile de seturi de date înregistrare:

- Tranzacții, date specifice cosurilor de cumpărare
 - exemplu: într-un magazin, setul de produse
 - se analizează asocierea între produsele în
 - posibilitate de reprezentare: indicator boolean
 - variație: câte exemplare din produs au fost

- Matrice de date:
 - pentru cazul în care datele au același set
 - fiecare dată în parte poate fi considerată o
 - fiecare atribut considerat este o dimensiune
 - este tipul de date standard pentru analiză
 - nota: între conceptul de dimensiune așa cum

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transformare
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Seturi
 - 1.1.4.2 Seturi
 - 1.1.4.3 Seturi
 - ▼ 1.2 Calitatea datelor
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantionare
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretizare
 - 1.3.7 Transformare
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de seturi
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecvență
 - 1.4.3.2 Procent
 - 1.4.3.3 Medie
 - 1.4.3.4 Măsură
 - 1.4.4 Statistică

**Proiect
of x L**

10.23

12.65

* Matrice de date rarefiate: * caz special a trasaturilor * exemplu: tranzactii din cosuri dintr-un anumit lexic; pentru un document sau nu = matricea document-termen * in (<https://docs.scipy.org/doc/scipy/reference>)

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transform
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Seturi
 - 1.1.4.2 Seturi
 - 1.1.4.3 Seturi
 - ▼ 1.2 Calitatea datelor
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transform
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

Docume
Docume
Docume

▼ 1.1.4.2 Seturi de date de tip graf

- reprezentare convenabila pentru cazurile:
 - graful reprezinta relatii între obiecte
 - obiectele însele sunt reprezentate ca graf

Caz 1: datele reprezinta relatii între obi

- obiectele sunt reprezentate ca noduri în graf
- relatiile dintre obiecte sunt reprezentate sub fo
- exemplu: pagini web care contin legaturi catre
- exemplu de algoritm ce foloseste structura de

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transform
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Seturi
 - 1.1.4.2 Seturi
 - 1.1.4.3 Seturi
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transform
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

Caz 2: obiectele-data sunt grafuri

- obiectele pot contine subobiecte care sunt leg
- uneori nu doar legaturile sunt importante, ci si
- exemplu: formulele chimice - benzen = C_6H_6
- utilitate: se poate detecta care substructura aș
- anumitor proprietati chimice.
- domeniu aparte: "minieritul" substructurilor

▼ 1.1.4.3 Seturi de tip secventa

- atributele au relatii care implica ordonare în tir
- subtipuri: date secventiale, secventa, serii de

Caz: date secventiale

- numite si date temporale
- fiecare înregistrare are un atribut suplimentar

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transform
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Seturi
 - 1.1.4.2 Seturi
 - 1.1.4.3 Seturi
 - ▼ 1.2 Calitatea datelor
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transform
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

Caz: date secventa

- setul de date entitati individuale, precum secv
- similare cu cele secventiale, dar fara timp incl
- pozitia din secventa este importanta
- exemplu: informatia genetica este o secventa
- aplicatie: predictia similaritatilor în structura si

▼ **Caz: serii de timp**

- fiecare înregistrare e o serie de timp = o serie
- exemplu: seturi de date de tip financiar, reprez
- exemplu: date meteo masurate lunar

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transforr
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setui
 - 1.1.4.3 Setui
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reduceri
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transforr
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

▼ **Caz: Date spatiale**

- cazul datelor care au attribute spatiale sau are
- exemplu: date climatice raportate pe regiuni
- exemplu: date adunate pentru scurgerea unui

▼ 1.2 Calitatea datelor

▼ 1.2.1 Probleme legate de masurarea

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setul
 - 1.1.4.3 Setul
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

- Presupunerea ca datele pe baza carora se fac
- Prevenirea problemelor care duc la scaderea
- Abordari:
 - detectarea si corectarea erorilor = curatar
 - construirea de algoritmi care sa tolereze c
- surse de probleme în calitatea datelor:
 - procesele de masurare
 - aplicatiile folosite
- Zgomotul
 - componenta aleatoare care se adauga ur
 - Exemplu: distorsiunea vocii unei persoan
 - Daca eroarea apare mereu în acelasi loc:

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

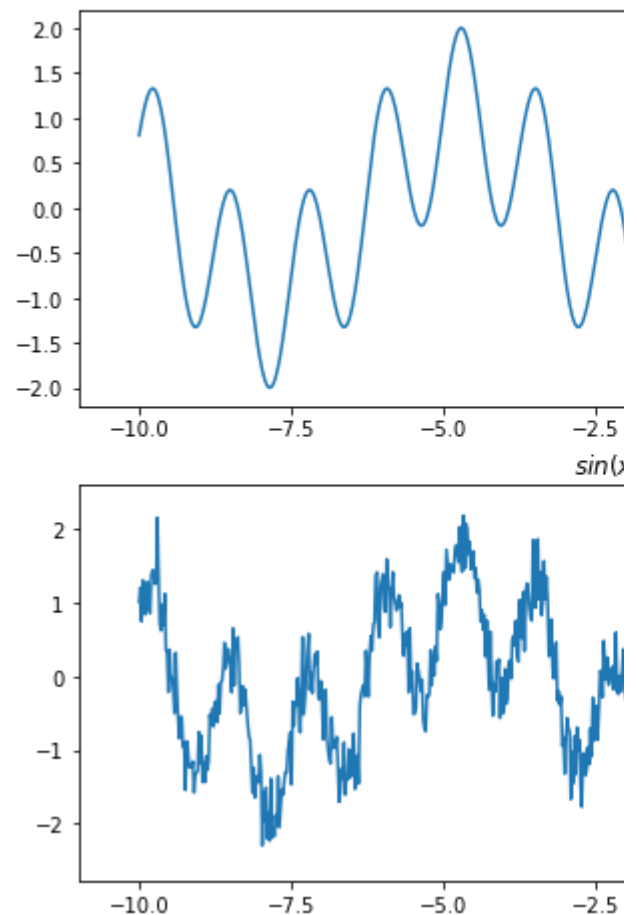
In [16]:

```
x = np.linspace(-10, 10, 1000)
y = np.sin(x) + np.sin(5*x)

fig, (ax1, ax2) = plt.subplots(2, figsize=(10, 10))
fig.suptitle('Date cu si fara zgomot')
ax1.set_title(r'$sin(x) + sin(5x)$')
ax1.plot(x, y)
ax2.set_title(r'$sin(x) + sin(5x) + 0.3 * np.random.randn(1000)$')
ax2.plot(x, y + 0.3 * np.random.randn(1000))
```

executed in 447ms, finished 17:33:07 2021-03-21

Da



1.2.2 Calitatea datelor: precizie, abatere

- **Precizie:** apropierea valorilor rezultate prin măsurători de la valoarea așteptată
- **Abatere (eng: bias):** o variație sistematică a măsurătorilor
- Exemplu: se măsoară o cantitate de 1 gram. \(\sigma = 0.001\).
- Precizia este considerată abaterea standard:
 - $\sigma = \sqrt{E[(X - E(X))^2]}$ deci pentru date
- **Acuratetea:** apropierea măsurătorilor față de valoarea așteptată

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transforr
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setui
 - 1.1.4.3 Setui
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transforr
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

▼ 1.2.3 Calitatea datelor: anomalii

- Anomaliile sunt obiecte cu caracteristici consi
- Anomaliile (outliers) nu sunt zgomote, ci obiec
- Utilitate: detectarea de nise pe piata, detectar

▼ 1.2.4 Calitatea datelor: valori lipsa

- Cazuri: una sau mai multe valori de attribute lip
- Motive pentru lipsa valorilor:
 - informatia nu este colectata — oamenii n
 - attributele nu se pot aplica tot timpul tutur
- Operarea în aceste situatii:
 - eliminarea obiectelor-data sau a atributel
 - estimarea valorilor lipsa
 - ignorarea valorilor lipsa în timpul analizei

▼ 1.2.5 Calitatea datelor: valori incons

- Valori inconsistente:
 - Exemplu: oras si cod postal precizate, da
 - Exemplu: typos, marimi cu valori impropri
 - Operare: detectarea valorilor gresite si co
 - E necesara utilizarea surselor de date rec
- Date duplicate:
 - Duplicarea poate sa fie exacta sau aproa
 - Exemplu: aceeași persoana cu adrese de
 - Procesul de curatare = deduplicare

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transform
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Seturi
 - 1.1.4.2 Seturi
 - 1.1.4.3 Seturi
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transform
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

1.2.6 Calitatea datelor din perspectiv

- Din perspectiva aplicatiilor, “datele au calitate
- Caracterul oportun al datelor — dacă datele s
- Relevanța — în cazul în care se vrea crearea
- alta situație este dată de esantionarea neade
- Cunoștințele apriori despre date — de exemplu
- reducerea redundanței și a dimensionalității; c

1.3 Preprocesarea datelor

- Strategii și tehnici complexe, ce pot cere până
- Două variante:
 1. selectarea obiectelor-dată și a atributelor
 2. crearea/schimbarea de atribute
- Variante de preprocesare:
 - agregare
 - esantionare
 - reducerea dimensionalității
 - selectarea unui subset de atribute
 - crearea de atribute
 - discretizare și binarizare
 - transformarea variabilelor

1.3.1 Agregare

- Scop: combinarea a două sau mai multe atri
- Utilitate:
 - reducerea cantității de date
 - schimbarea scalei: orasele sunt agregate
 - date mai stabile: datele agregate au tendi

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transforr
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setui
 - 1.1.4.3 Setui
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transforr
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

▼ 1.3.2 Esantionare

- principala tehnica folosita pentru selectarea d:
- în statistica, a fost folosita atât pentru investig
- este folosita pentru ca obtinerea întregului set
- în DM esantionarea este folosita pentru ca prc
- un esantion este reprezentativ daca are aprox
- utilizarea unui esantion reprezentativ e aproaș
- tipuri de esantionare:

Contents 🔄 ⚙️

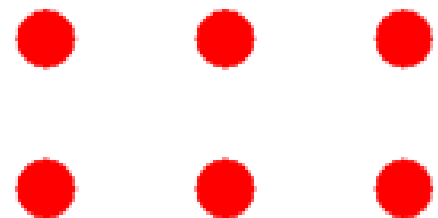
- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transform
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Seturi
 - 1.1.4.2 Seturi
 - 1.1.4.3 Seturi
 - ▼ 1.2 Calitatea datelor
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transform
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

- esantionare aleatoare uniforma: avem o p
- esantionare fara înlocuire: daca un obiect
- esantionare cu înlocuire: obiectele nu sur
- esantionarea stratificata: se divid datele în
- cu proportia din multimea initiala.



Esantion de 8000 de puncte

- Determinarea marimii esantionului
 - problema: avem un set de date care cons
 - trebuie extras un esantion astfel încât din
 - “garantarea” se exprima probabilist: care
 - anterior sa depaseasca un anumit prag?



▼ 1.3.3 Reducerea dimensionalitatii

- seturile de date pot avea un numar mare de a
- valorile pentru densitate de probabilitate si dis

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transform
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Seturi
 - 1.1.4.2 Seturi
 - 1.1.4.3 Seturi
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transform
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

- soluție: reducerea dimensionalității fără a pier
- beneficii: algoritmi de DM lucrează mai eficient
- Scop:
 - evitarea blestemului dimensionalității
 - reducerea timpului de rulare necesar alg
 - datele devin mai ușor de vizualizat
 - poate ajuta la eliminarea trăsăturilor irelev
- Tehnici folosite:
 - analiza componentelor principale (Princip
 - descompunerea valorilor principale (Singl
 - alte metode: supervizate și transformări n

▼ 1.3.4 Selectarea subsetului de atribu

- motivatie: de multe ori, în seturile de date pot
- unele din aceste atribute pot fi eliminate prin “
- varianta ideală de lucru: se încearcă toate cor
- variante:
 1. metode încorporate — algoritmul de DM î
 2. metode de filtrare — atributele sunt selec
 3. metode bazate pe încercare — se foloseș
- Pentru selectarea subsetului de trăsături e ne
 - o măsură pentru evaluarea unui subset d
 - o metodă de căutare pentru generarea ur
 - un criteriu de oprire
 - o procedură de validare
- Alternativa pentru selectarea de atribute: ponc
 - atribute importante → ponderi mari; atribu
 - ponderile se pot asocia pe baza cunoștinț

▼ 1.3.5 Crearea de trăsături

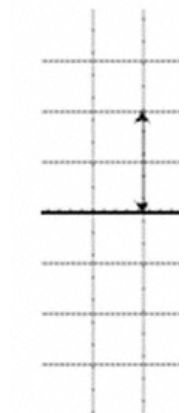
- Scop: crearea de noi atribute pe baza celor ex
- Noile atribute pot releva mai eficient informaț
- Numărul atributelor noi obținute poate fi mai m
- Metode:
 - Extragere de trăsături
 - Transformarea datelor, utilizarea unui alt :
 - Construirea de trăsături
- Exemplu: extragerea de trăsături
 - crearea unui nou set de atribute pe baza
 - exemplu: fotografie → detectarea muchiilor
 - metodele sunt strâns legate de domeniu -

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setul
 - 1.1.4.3 Setul
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic



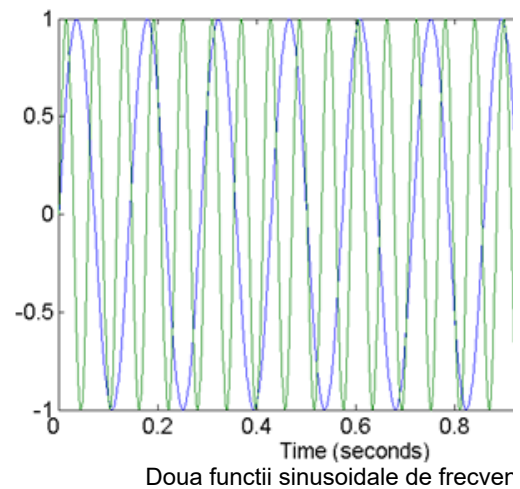
- Exemplu: transformarea datelor
 - trecerea de la coordonatele carteziene la



- transformata Fourier

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic



- construirea de noi atribute
 - uneori, atributele originare pot sa aib
 - exemplu: pentru retele neurale intrari
 - exemplu de trasatura derivata: densi
 - crearea de noi trasaturi necesita cun

Discretizare si binarizare

```

* unii algoritmi cer ca datele sa fie în
* algoritmi de determinare a asociierilor
* apare nevoia de a transforma atribute c
* Discutie pentru binarizare:
    * $m$ valori -> fiecare valoare initi
transcrierea în baza 2
    * daca atributele sunt ordinale, atur
    * exemplu: {slab, mediu, bun} -> ${\{
    * $\lceil \log_2 m \rceil$ biti folos
    * one hot encoding: pentru m clase se
    exemplu: {slab, mediu, bun} -> $(x_1
* Discretizare a atributelor continue
    * procesul de discretizare este deper
    * pasi:
        * se decide numarul de categorii
        * se decide modul de asociere înt
        * exemplu: se sorteaza datele, se
eticheta
* Discretizare nesupervizata
    * intervalele pot fi de latime egala
    * prezinta frecvente egale
    * metode de clustering, e.g. K-means
    * inspectare vizuala a datelor
<table>
<tr>
<td> 
<tr>
<td> <img src="./images/discretizar
<td> <img src="./images/discretizar

```


Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

```

</tr>
</table>
* Discretizare supervizata
  * utilizarea de informatie suplimenta
  * conceptual: discretizarea sa se fac
  * metode statistice: intervale mai mi
* Atribute categoriale cu prea multe atri
  * exemplu: o multitudine de departame
  * se pot unifica pe domenii: ingineri
  * altfel: gruparea poate fi facuta în
clasificare

```

1.3.7 Transformarea de variabile

- transformare care se aplica la toate valorile ur
- exemplu: daca doar amplitudinea unei valori e
- functii folosite: $\exp(x)$, $\log(x)$, \sqrt{x} , $1/x$, func
- functii complexe: normalizare
- pentru datele care prezinta domeniu mare de
 - exemplu: transferuri de date de 109, 108, 9, 8, 3, 1.
- aplicarea trebuie facuta în cunostinta de cauz

1.3.8 Exemplu de transformare: stan

- în statistica: daca \bar{x} este media unui atribut si
- atunci transformarea $x \rightarrow \frac{x - \bar{x}}{s_x}$ creeaza o noi
- exemplu: persoane cu variabilele: venitul anual
- venituri ar domina distanta între doi oameni, d
- problema: abaterea patratica medie e influent

unde μ e media sau mediana

Explorarea datelor: Statisti

1.4.1 Explorarea datelor

- Explorarea datelor reprezintă investigarea pre
- Pasul de explorare poate fi de folos în alegere
- Se poate folosi abilitatea naturală a oamenilor
- Domeniul a fost introdus de către statisticianu
- AED este domeniu opus lui "Confirmatory Dat

- Curs de AED: [aici \(https://www.coursera.org/learn/science&utm_source=gg&utm_medium=sem&ROW&adgroupid=116274867101&device=c&lyWJrboZGvrVzGL9QL0-8PazMsYrrQaAnG4E](https://www.coursera.org/learn/science&utm_source=gg&utm_medium=sem&ROW&adgroupid=116274867101&device=c&lyWJrboZGvrVzGL9QL0-8PazMsYrrQaAnG4E)

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

- * În AED, așa cum este definit de Tukey:
 - * Focus-ul este pe vizualizare
 - * Gruparea (clustering) și detectarea
 - * Acestea două sunt subdomenii aparte

Setul de date Iris

- * Constă în date măsurate pentru 150 de f
- * Măsurătorile sunt pentru lungimea/lățim
- * A cincea coloană este specia florii - a
- * Datele se pot descărca [[de aici](#)](http://

In [35]:

```
iris_url = 'http://archive.ics.uci.edu/iris_data = pd.read_csv(iris_url, header = 0) data.columns = iris_columns data.describe(include='all')
```

executed in 724ms, finished 20:37:59 2021-03-21

Out[35]:

	sepal_length	sepal_width	petal_length
count	150.000000	150.000000	150.000000
unique	NaN	NaN	NaN
top	NaN	NaN	NaN
freq	NaN	NaN	NaN
mean	5.843333	3.054000	3.758667
std	0.828066	0.433594	1.764420
min	4.300000	2.000000	1.000000
25%	5.100000	2.800000	1.600000
50%	5.800000	3.000000	4.350000
75%	6.400000	3.300000	5.100000
max	7.900000	4.400000	6.900000

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transforr
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setui
 - 1.1.4.3 Setui
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transforr
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

iris setosa



petal sepal

Sursa: morioh.com

1.4.3 Statistici de sumarizare

- Statisticile de sumarizare sunt numere care sc
- Reprezintă manifestarea cea mai vizibilă a sta
- Exemple: frecvența, media, dispersia

1.4.3.1 Frecvența și valoarea modală

Pentru un set de m date categoriale cu valorile $\{v$

iar frecvența absolută este numărul fracției de

In [30]:

```
from collections import Counter
freq_iris = Counter(data['class'])
rel_freq_iris = {iris_type: freq / len(d
```

executed in 5ms, finished 20:28:35 2021-03-21

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

In [31]:

```
print(f'Frequency: {freq_iris}')
print(f'Relative Frequency: {rel_freq_i
```

executed in 6ms, finished 20:29:02 2021-03-21

Frequency: Counter({'Iris-setosa': 50, 'Iris-versicolour': 49, 'Iris-virginica': 51})
 Relative Frequency: {'Iris-setosa': 0.3333333333333333, 'Iris-versicolour': 0.32666666666666666, 'Iris-virginica': 0.34}

* Valoarea modală (sau moda) este valoarea care apare cel mai des în datele
 \$\$\$\$
 moda = \arg\max\limits_{v_i} frecventa(v_i)
 \$\$\$\$
 * Atenție la situația când o anumită valoare apare de mai multe ori
 * Pot exista seturi de date pentru care nu există o valoare modală
 * Pentru valori continue, conceptele de modă și medie nu sunt aplicabile

Percentile

* Pentru cazul valorilor ordonate se pot calcula percentilii
 * Pentru un atribut continuu sau ordinal se poate calcula percentila
 șirul de valori ale lui \$x\$ astfel încât \$p\%\$ din valorile sunt mai mici sau egale cu \$x\$
 * Nu există o definiție standardizată pentru calculul percentilelor
 * Pentru cazul în care se calculează percentilele pentru un set de date, acestea devin neesențiale
 * Tradițional se consideră \$x_{(0\%)} = \min(x)\$ și \$x_{(100\%)} = \max(x)\$
 * Mod de calcul pentru determinarea celei mai apropiate de \$\frac{m}{100}\$ din sortat
 * Funcție utilizabilă: [\[numpy.percentile\]](#) sau [\[pandas.DataFrame.quantile\]](#) (<https://pandas.pydata.org/pandas-docs/stable/10min.html#percentiles>)

In [34]:

```
help(np.percentile)
```

executed in 182ms, finished 20:36:59 2021-03-21

```
a = np.arange(4)
p = np.linspace(0, 100, 6001)
ax = plt.gca()
lines = [
    ('linear', None),
    ('higher', '--'),
    ('lower', '--'),
    ('nearest', '-.'),
    ('midpoint', '-.'),
]
for interpolation, style in lines:
    ax.plot(
        p, np.percentile(a, p, interpolation),
        label=interpolation, linestyle=style
    )
ax.set(
    title='Interpolation methods',
    xlabel='Percentile',
    ylabel='List item returned',
    yticks=a
)
```

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

In [53]:

```
_25th_percentile_sepal_length = np.perc

print(f'Cate valori sunt strict mai mic
print(f'Cate valori sunt strict mai mic
print(f'lista valorilor sortate: {list(

# se poate intelege pentru aceste date
```

executed in 9ms, finished 20:44:19 2021-03-21

Cate valori sunt strict mai mici decat a
 Cate valori sunt strict mai mici sau egal
 lista valorilor sortate: [4.3, 4.4, 4.4,
 9, 4.9, 5.0, 5.0, 5.0, 5.0, 5.0, 5.0, 5.0,
 5.4, 5.4, 5.4, 5.4, 5.4, 5.4, 5.5, 5.5, 5.
 5.7, 5.8, 5.8, 5.8, 5.8, 5.8, 5.8, 5.8, 5.
 6.2, 6.3, 6.3, 6.3, 6.3, 6.3, 6.3, 6.3, 6.
 6.7, 6.7, 6.7, 6.7, 6.7, 6.7, 6.8, 6.8, 6.

In [54]:

```
data['sepal_length'].quantile(q=0.25)
```

executed in 15ms, finished 20:48:41 2021-03-21

Out[54]: 5.1

Media si mediana

* Pentru un set de valori $\{x_1, x_2, \dots, x_n\}$

$$\overline{x} = E(x) = \text{mean}(x) = \frac{1}{n} \sum_{i=1}^n x_i$$

 * Mediana este a 50-a percentila. O variabilă ordonată (reordonarea) $\{x_{(1)}, x_{(2)}, \dots, x_{(n)}\}$

$$\text{mediana}(x) = \begin{cases} x_{(r)} & \text{dacă } n \text{ este par} \\ \frac{x_{(r)} + x_{(r+1)}}{2} & \text{dacă } n \text{ este impar} \end{cases}$$

 O varianta mai rafinata (si care se poate calcula in timp complexitate liniara in loc de $O(n \log n)$)
 Obtinere: [\[numpy.median\]](https://numpy.org/doc/stable/reference/generated/numpy.median.html) (<https://numpy.org/doc/stable/reference/generated/numpy.median.html>)
https://pandas.pydata.org/pandas-docs/stable/10min/05_datamanipulation.html#pandas.Series.quantile

In [57]:

```
np.median(data.sepal_length.values)
```

executed in 13ms, finished 21:18:38 2021-03-21

Out[57]: 5.8

Diferente intre medie si mediana:

- Media este valoare "de mijloc" doar dacă distribuția este simetrică
- Dacă distribuția este asimetrică, atunci mediana este o măsură mai bună a valorii centrale

- Media este influențată de outliers, în timp ce n
- Medie retezată (eng: trimmed mean) se utilize din date; se calculează media pentru ceea ce
- media standard se obține din media retezată c

Contents 🔄 ⚙️

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecv
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

Exemple:

```
* Considerăm valorile {1, 2, 3, 4, 5, 90}
considerabil diferită față de media setul
* Media, medianele și valoarea de trimmec
<img src="./images/stats_iris.png" alt="c
* Exercițiu: dacă valoarea medianei este
```

Măsurari ale împrăstierii datelor

```
* Sunt măsuri care cuantifică concentrare
* Diametrul domeniului de valori (eng: range)
$$
range(x) = \max(x) - \min(x) = x_{(m)} - x_{(1)}
$$
* Range-ul este nerelevant, deoarece pute
outlier măresc artificial diametrul setul
* Varianța (dispersia; eng: variance) un
$$
varianța(x) = s_x^2 = \frac{1}{m-1} \sum (x_i - \bar{x})^2
$$
* Utilizarea numitorului $m - 1$ în loc c
și are ca scop corectarea abaterii din es
* Abaterea standard este $s_x = \sqrt{s_x^2}$
* Python: \[numpy.std\]\(https://numpy.org/c
```

In [59]: `np.std(data.petal_length), np.std(data.`

executed in 17ms, finished 21:40:36 2021-03-21

Out[59]: (1.7585291834055201, 0.8253012917851409)

- Deoarece media poate să fie distorsionată de
- Se preferă considerarea altor trei măsuri:
 - absolute average deviation, AAD:
 - median absolute deviation, MAD:
 - interquartile range:

Verificam ca prin standardizare de date,

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic

In [69]:

```
med_petal_length = np.mean(data.petal_l
# centram datele: media lor devine 0
centered_petal_length = data.petal_leng
# assert 0 == np.mean(centered_petal_le
assert np.allclose(0, np.mean(centered_
std_petal_length = np.std(data.petal_le
std_centered_petal_length = np.std(cent
assert np.allclose(std_petal_length, st
standardized_petal_length = centered_pe
# verificare: media 0...
assert np.allclose(0, np.mean(standardi
# si standard deviation ~ 1
assert np.allclose(1, np.std(standardiz
```

executed in 52ms, finished 21:53:56 2021-03-21

1.4.4 Statistici de sumarizare a datel

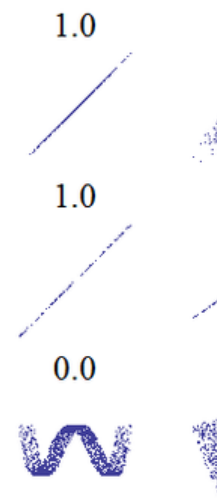
- Date multivariate: date cu mai multe atribute
- Pentru atributul x_i calculăm media \bar{x}_i
- Media setului de obiecte este $\mathbf{x} = (\bar{x}_1, \dots, \bar{x}_n)$
- Analog se poate calcula dispersia, mediana ei
- Matricea de covarianță: elementul s_{ij} este cov

unde x_{pq} este a p -a valoare a atributului x_q

- s_{ij} este măsură a gradului în care două atribu
- $s_{ij} = 0$ arata ca atributele s_i si s_j nu sunt lini
- pentru a elimina dependenta de magnitudinea

- r_{ij} se mai numeste corelatia Pearson a atribu
- $r_{ij} = \pm 1$ indica faptul ca x_i e in relatie liniara

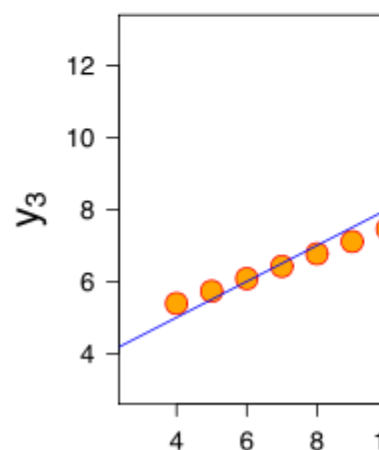
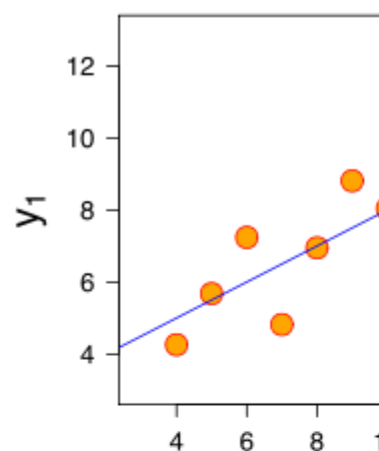
cu $\text{sgn}(a) = \text{sgn}(r_{ij})$



Seturi de date (x, y) împreună cu coeficientul de corelație (ρ) (sus) precum și modul în care ele sunt legate liniar relații mai complexe între date (rândul de jos). Sur:

Legat de coeficientul de corelație, câteva observații:

- “Corelația nu înseamnă cauzalitate” – nu se poate spune că o corelație mare poate fi o condiție necesară pentru a punct de pornire în cercetarea unei legături între date.
- Corelația și liniaritatea – coeficientul Pearson
- Exemplu: 4 seturi de date cu două atribute; în fiecare caz (0.816); cu toate acestea, legătura



Date cu caracteristici numerice identice (medi

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transformare
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Seturi
 - 1.1.4.2 Seturi
 - 1.1.4.3 Seturi
 - ▼ 1.2 Calitatea datelor
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Eliminare
 - 1.3.3 Reducere
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretizare
 - 1.3.7 Transformare
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de descriere
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frecvență
 - 1.4.3.2 Percentil
 - 1.4.3.3 Medie
 - 1.4.3.4 Măsură
 - 1.4.4 Statistică

Contents

- ▼ 1 Curs 5. Tipuri de
 - ▼ 1.1 Tipuri de date
 - 1.1.1 Tipuri de
 - 1.1.2 Transfor
 - 1.1.3 Descriere
 - ▼ 1.1.4 Tipuri de
 - 1.1.4.1 Setur
 - 1.1.4.2 Setu
 - 1.1.4.3 Setu
 - ▼ 1.2 Calitatea dat
 - 1.2.1 Probleme
 - 1.2.2 Calitatea
 - 1.2.3 Calitatea
 - 1.2.4 Calitatea
 - 1.2.5 Calitatea
 - 1.2.6 Calitatea
 - ▼ 1.3 Preprocesare
 - 1.3.1 Agregare
 - 1.3.2 Esantion
 - 1.3.3 Reducer
 - 1.3.4 Selectare
 - 1.3.5 Crearea
 - 1.3.6 Discretiz
 - 1.3.7 Transfor
 - 1.3.8 Exemplu
 - ▼ 1.4 Statistici de s
 - 1.4.1 Explorare
 - 1.4.2 Setul de
 - ▼ 1.4.3 Statistici
 - 1.4.3.1 Frec
 - 1.4.3.2 Perc
 - 1.4.3.3 Medi
 - 1.4.3.4 Măsi
 - 1.4.4 Statistic