



Facultatea de Electronică,
Telecomunicații și
Tehnologia Informației

SISTEME INTELIGENTE DE SUPORT DECIZIONAL

Ș.l.dr.ing. Laura-Nicoleta IVANCIU

Curs 6 – Clasificare. Machine Learning pentru sisteme de decizie.

Cuprins

- Clasificare – definire, tipuri
- Etape ale clasificării
- Machine Learning – definire, algoritmi, utilizare, limitări

Previously on SISD (C4)

Predictive Data Mining (P)

- Clasificare
- Regresie
- Serii temporale

Descriptive Data Mining (D)

- Detectie de anomalii
- Clustering
- Sumarizare (summarizing)
- Association rule learning

Previously on SISD (C4)

Clasificare (P)

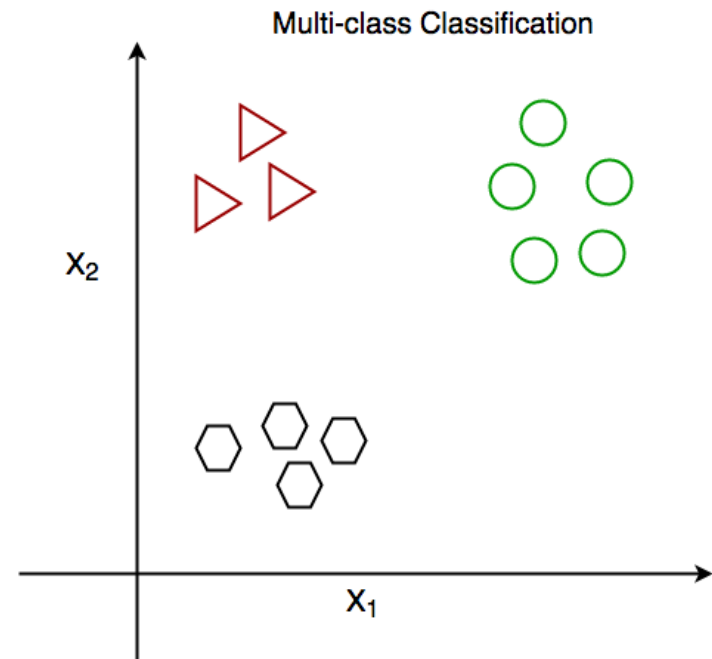
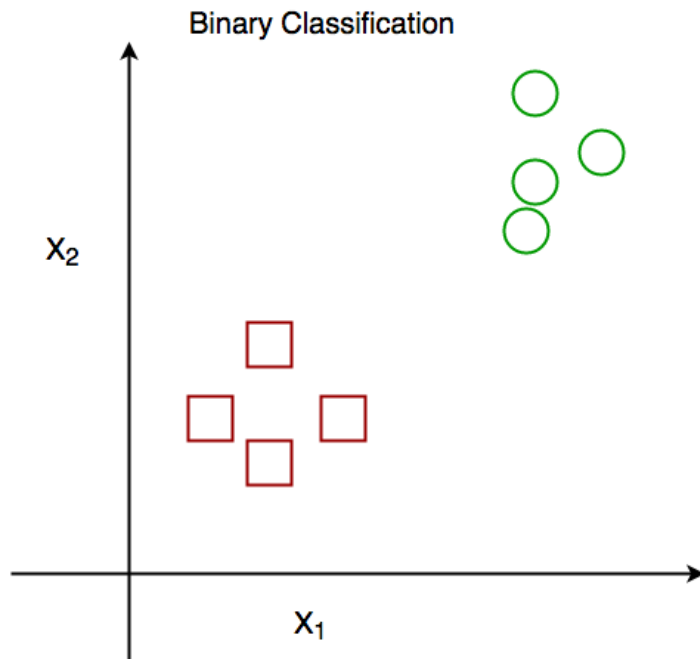
- binară/non-binară (multi-clasă)

Ex. spam/non-spam, mașină/camion/motocicletă/bicicletă

- împărțire pe categorii/clase, pe baza unor observații, proprietăți, trăsături (*explanatory variables, features*)
- algoritm de clasificare = clasificator

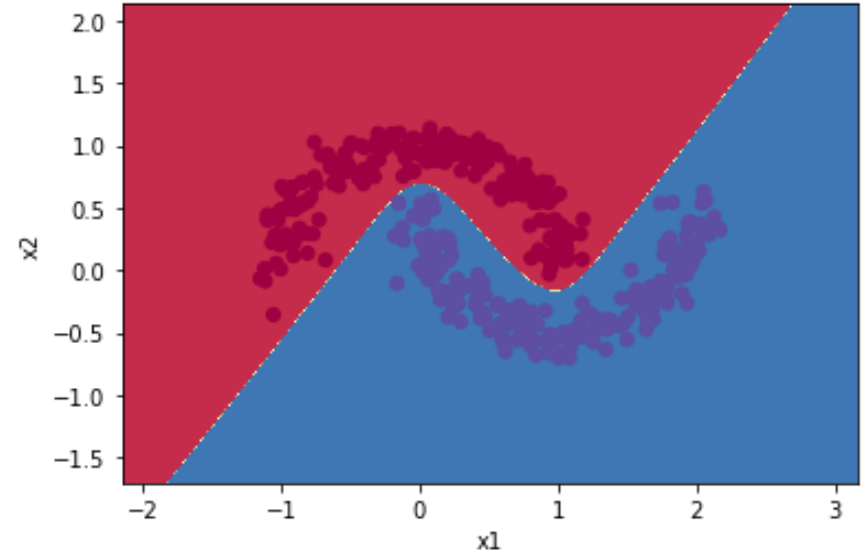
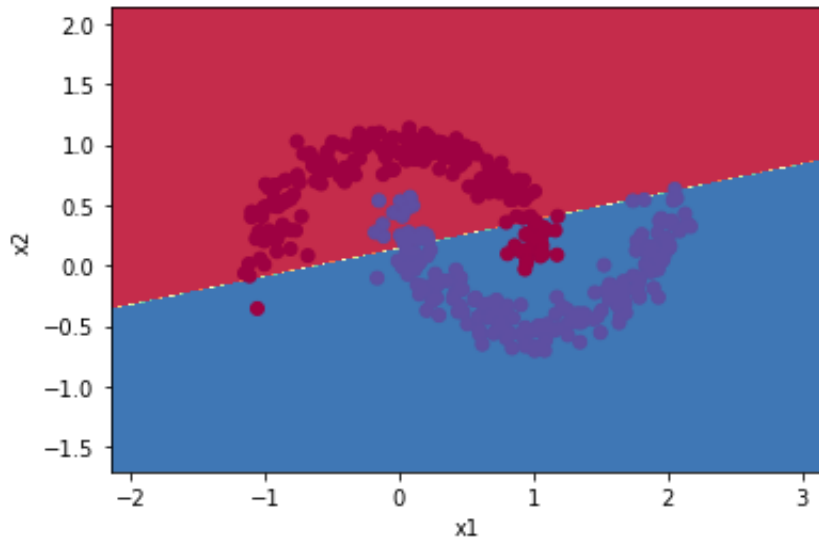
Previously on SISD (C4)

Clasificare (P)

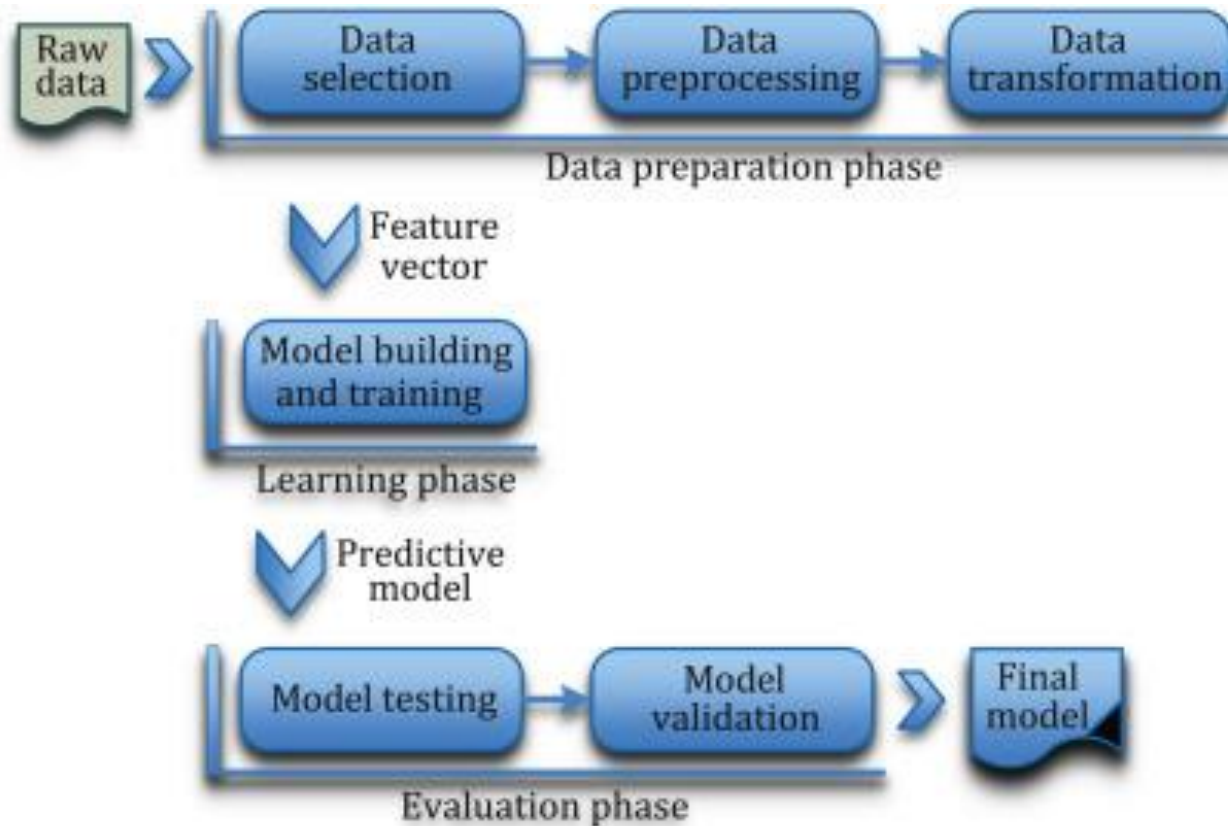
[Sursă](#)

Previously on SISD (C4)

Clasificare liniară/nelinară

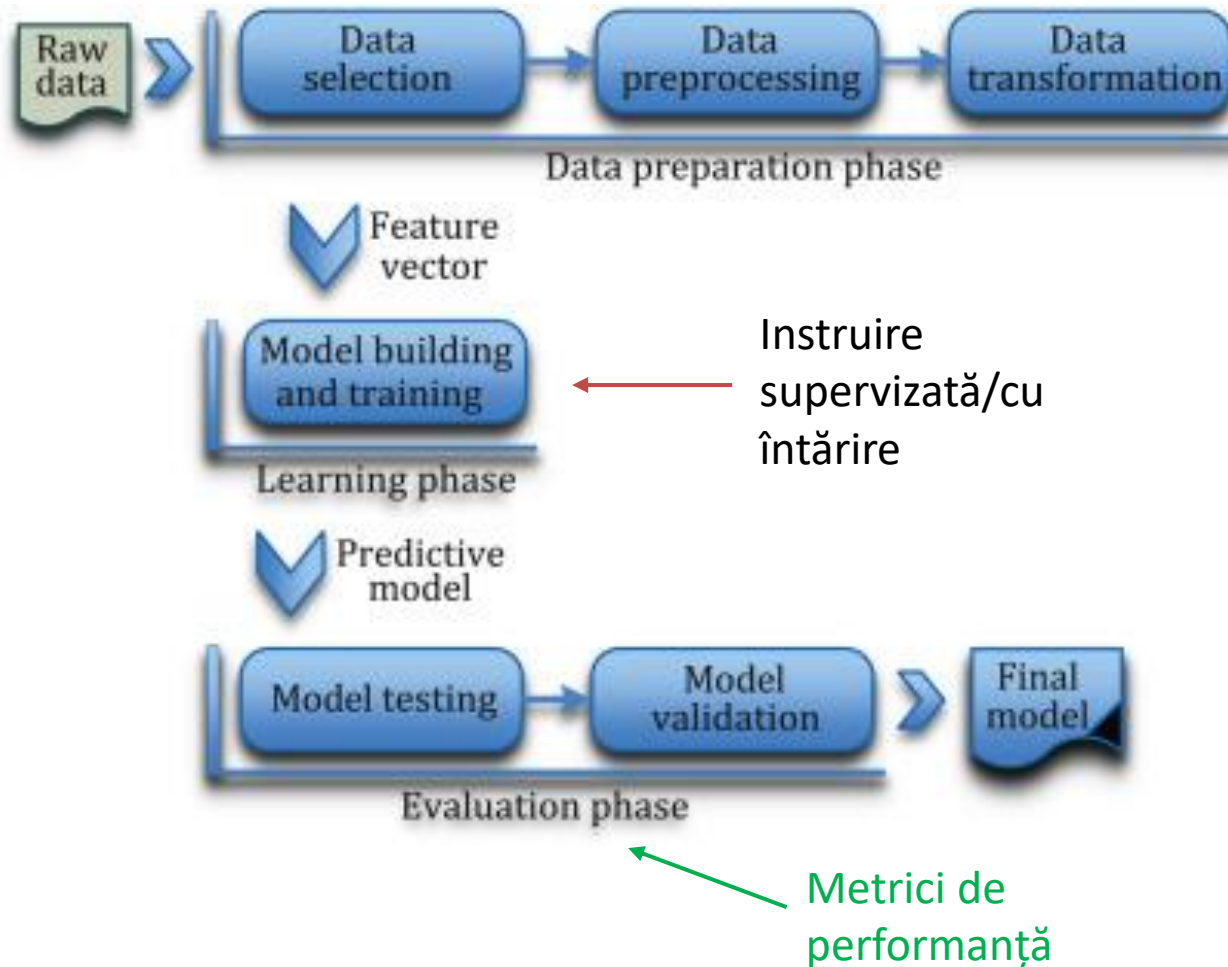


Etapе ale clasificării



[Sursă](#)

Etapе ale clasificării



[Sursă](#)

Etape ale clasificării

➤ Pregătirea datelor

- poate implica tehnici de Data Mining de tip descriptiv (sumarizare, detecție de anomalii, association rule learning)
- setul de date final este pregătit pentru utilizare la instruire, prin divizare în subseturi

Posibile variante de subseturi:

80% Train + 20% Test



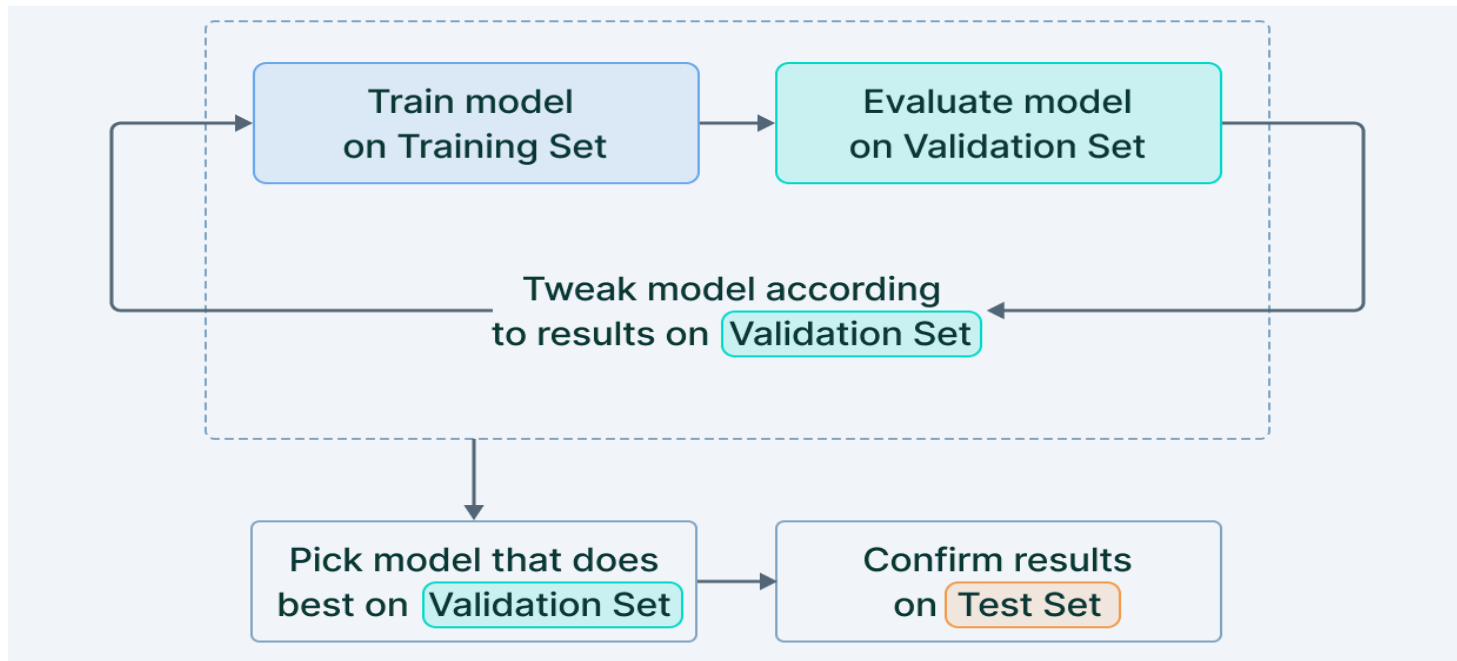
70% Train + 15% Valid + 15% Test



Etape ale clasificării

➤ Pregătirea datelor

Subseturile de antrenare/validare/testare



Etape ale clasificării

➤ Pregătirea datelor

Subseturile de antrenare/validare/testare

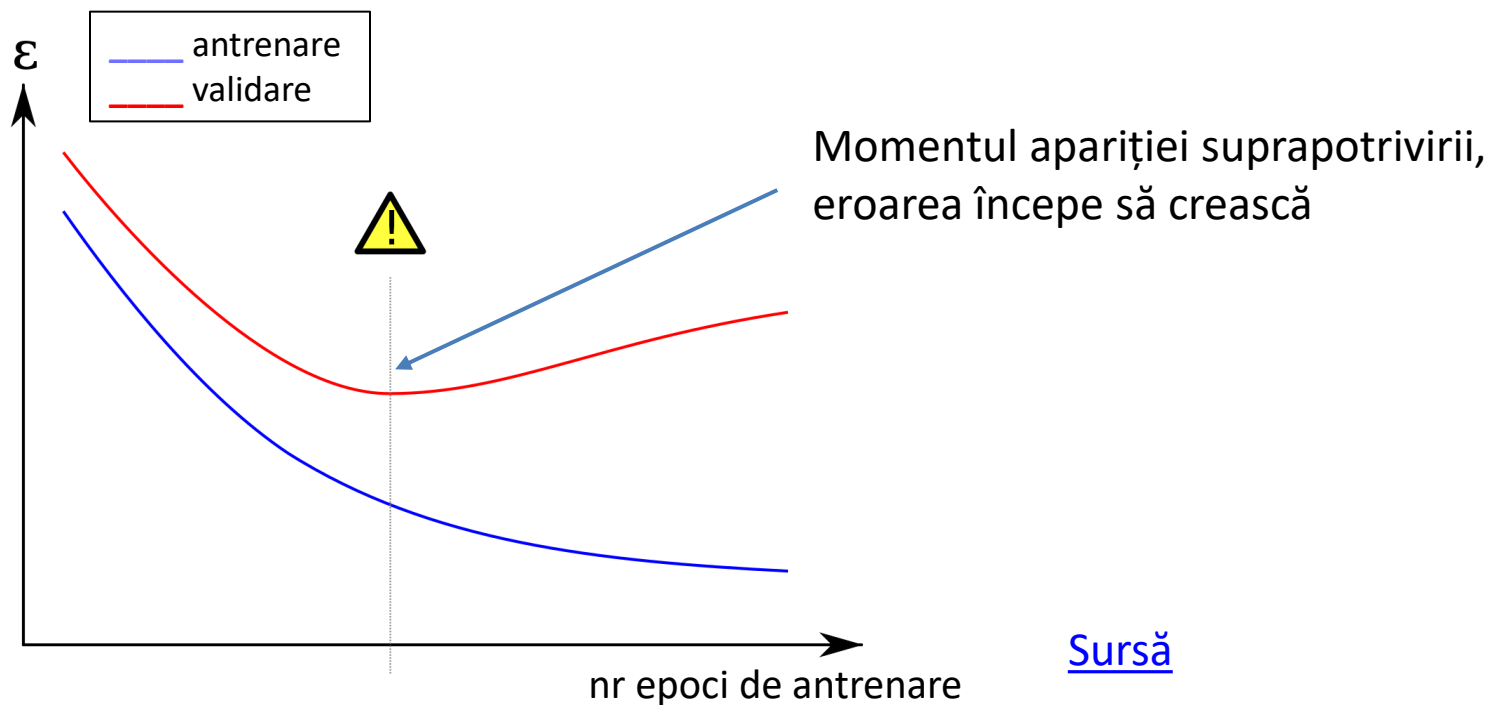
Observații:

- datele de validare + testare nu sunt incluse în procesul de instruire (modelul nu le “vede” pe durata instruirii)
- subsetul de validare previne apariția suprapotrivirii (overfitting)
- subsetul de testare evaluează performanța modelului

Etape ale clasificării

➤ Pregătirea datelor

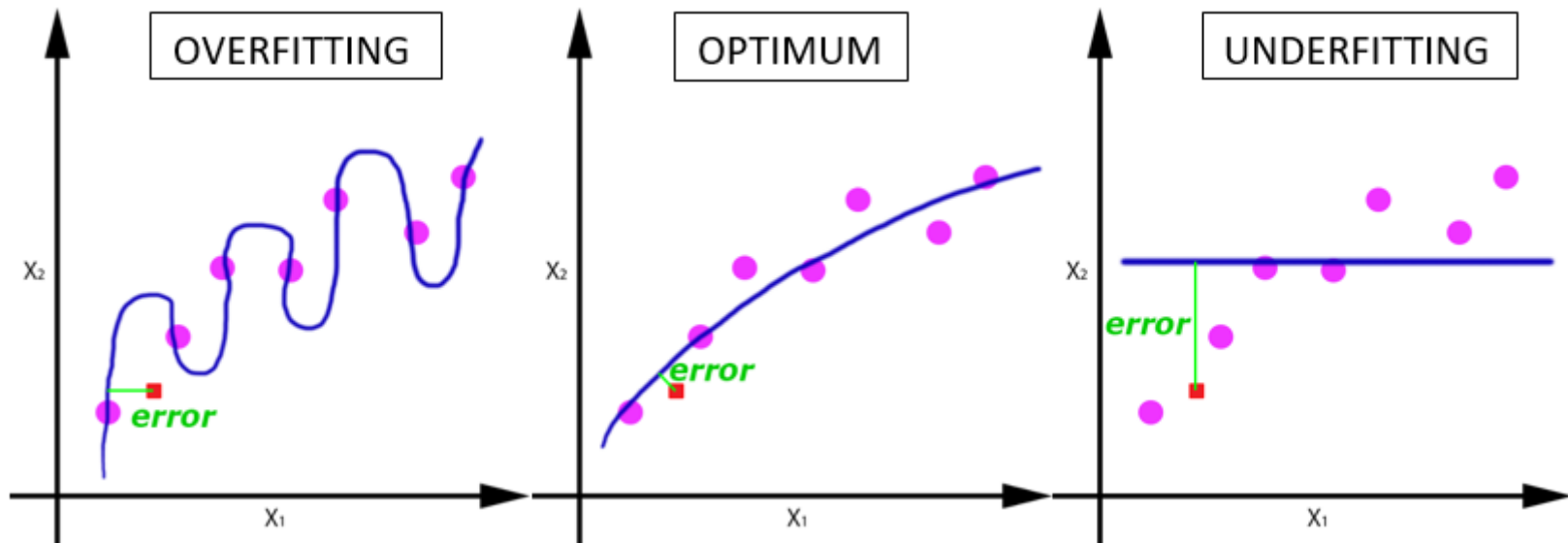
Subseturile de antrenare/validare/testare



Etape ale clasificării

➤ Pregătirea datelor

Subseturile de antrenare/validare/testare



Sursă

Etape ale clasificării

➤ Instruirea modelului

- epoci de instruire

Condiții de oprire: număr finit de epoci, eroare mai mică decât o valoare predefinită, număr de epoci consecutive fără îmbunătățiri

Epocă/iterație/batch

Epocă de instruire: întreg setul de date este trecut prin model

Batch (lot): o parte a setului de date

Batch size: dimensiune batch; determină frecvența de actualizare a parametrilor modelului

Iterație: un batch este trecut prin model

Etape ale clasificării

➤ Instruirea modelului

Epocă/iterație/batch

*Dimensiunea setului de date = nr. batch * batch_size*

Nr. de iterații necesare pentru a parcurge întreg setul de date = nr. batch

Nr. epoci = nr. iterații când batch_size = dimensiunea setului de date (nr. batch = 1)

Ex_1. Fie un set de date de 3000 de înregistrări.

Dacă *batch_size* = 500, de câte iterații este nevoie pentru o epocă de instruire?

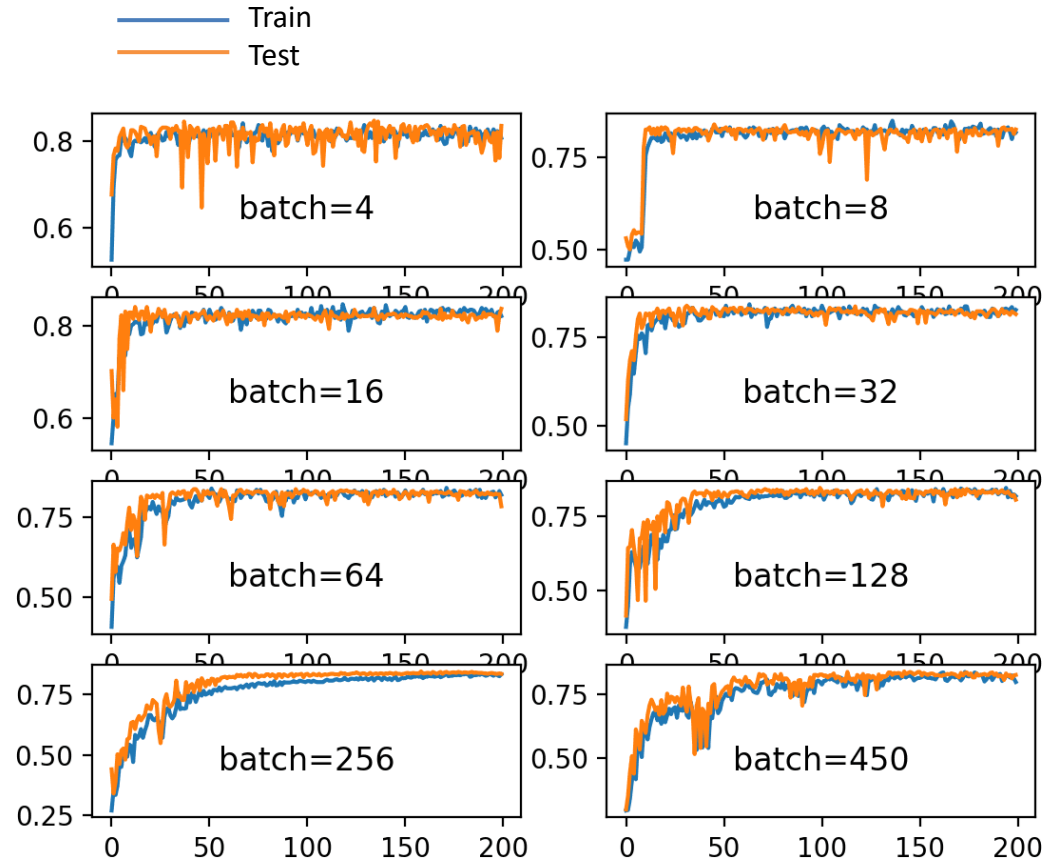
Ex_2. Dacă *batch_size* = 50 și este nevoie de 12 iterații/epocă, care este dimensiunea setului de date?

Etape ale clasificării

➤ Instruirea modelului

Epocă/iterație/batch

- Optimum batch size?
- De ce sunt necesare mai multe epoci de instruire?



Classification Accuracy on Train and Test Datasets With Different Batch Sizes [Sursă](#)

Etape ale clasificării

➤ Evaluarea performanțelor modelului

- pe setul de date de test (care nu a fost inclus în instruire)

Pe baza matricii de confuzie, se determină:

TP, TN, FP, FN, accuracy, precision, recall, F1-score

	Valori cunoscute (ground truth)				
		Clasa 1	Clasa 2	...	Clasa n
Valori estimate (predicted)	Clasa 1	5	2	...	1
	Clasa 2	1	7	...	1

	Clasa n	0	1	...	4

Etape ale clasificării

➤ Evaluarea performanțelor modelului

TP = adevăr pozitiv – un exemplar care aparține clasei a fost detectat ca aparținând clasei

- se calculează pentru fiecare clasă

TP pentru clasa 1 = 5

TP pentru clasa 2 = 7

	Valori cunoscute (ground truth)				
	Clasa 1	Clasa 2	...	Clasa n	
Valori estimate (predicted)	Clasa 1	5	2	...	1
	Clasa 2	1	7	...	1

	Clasa n	0	1	...	4

Etape ale clasificării

➤ Evaluarea performanțelor modelului

TN = adevăr negativ – un exemplar care nu aparține clasei nu a fost detectat ca aparținând clasei

- se calculează pentru fiecare clasă

TN pentru clasa 1 = 7

TN pentru clasa 2 = 5

TN este echivalent cu clasificarea corectă pentru clasa opusă

	Valori cunoscute (ground truth)				
	Clasa 1	Clasa 2	...	Clasa n	
Valori estimate (predicted)	Clasa 1	5	2	...	1
	Clasa 2	1	7	...	1

	Clasa n	0	1	...	4

Etape ale clasificării

➤ Evaluarea performanțelor modelului

FP = fals pozitiv – un exemplar care nu aparține clasei a fost detectat ca aparținând clasei

- se calculează pentru fiecare clasă

FP pentru clasa 1 = 2

FP pentru clasa 2 = 1

	Valori cunoscute (ground truth)				
	Clasa 1	Clasa 2	...	Clasa n	
Valori estimate (predicted)	Clasa 1	5	2	...	1
	Clasa 2	1	7	...	1

	Clasa n	0	1	...	4

Etape ale clasificării

➤ Evaluarea performanțelor modelului

FN = fals negativ – un exemplar care aparține clasei nu a fost detectat ca aparținând clasei (non-detectie)

- se calculează pentru fiecare clasă

FN pentru clasa 1 = 1

FN pentru clasa 2 = 2

	Valori cunoscute (ground truth)				
	Clasa 1	Clasa 2	...	Clasa n	
Valori estimate (predicted)	Clasa 1	5	2	...	1
	Clasa 2	1	7	...	1

	Clasa n	0	1	...	4

Etape ale clasificării

➤ Evaluarea performanțelor modelului

Accuracy = măsura clasificărilor corecte

$$\text{Accuracy} = \frac{\text{nr. clasificări corecte}}{\text{nr. total de date}}$$

	Valori cunoscute (ground truth)				
	Clasa 1	Clasa 2	...	Clasa n	
Valori estimate (predicted)	Clasa 1	5	2	...	1
	Clasa 2	1	7	...	1

	Clasa n	0	1	...	4

Etape ale clasificării

➤ Evaluarea performanțelor modelului

Precision = măsura falselor clasificări

$$Precision = \frac{TP}{TP + FP}$$

FP = 0 rezultă $Precision = 100\%$

Recall = măsura non-deteecțiilor

$$Recall = \frac{TP}{TP + FN}$$

FN = 0 rezultă $Recall = 100\%$

	Valori cunoscute (ground truth)				
	Clasa 1	Clasa 2	...	Clasa n	
Valori estimate (predicted)	Clasa 1	5	2	...	1
	Clasa 2	1	7	...	1

	Clasa n	0	1	...	4

Etape ale clasificării

➤ Evaluarea performanțelor modelului

F1 score

$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall}$$

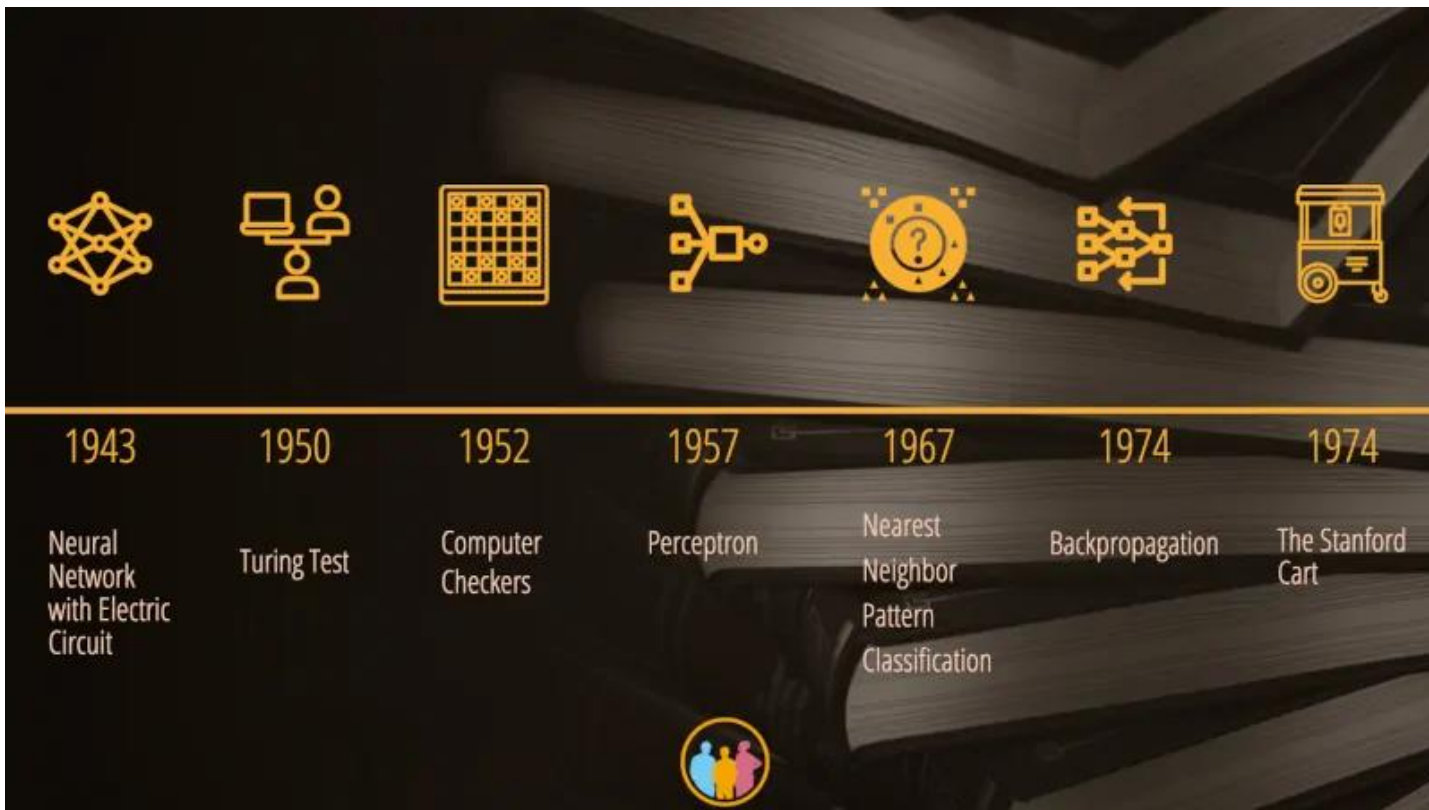
Care metrică este cea mai potrivită?

	Valori cunoscute (ground truth)				
		Clasa 1	Clasa 2	...	Clasa n
Valori estimate (predicted)	Clasa 1	5	2	...	1
	Clasa 2	1	7	...	1

	Clasa n	0	1	...	4

Machine Learning

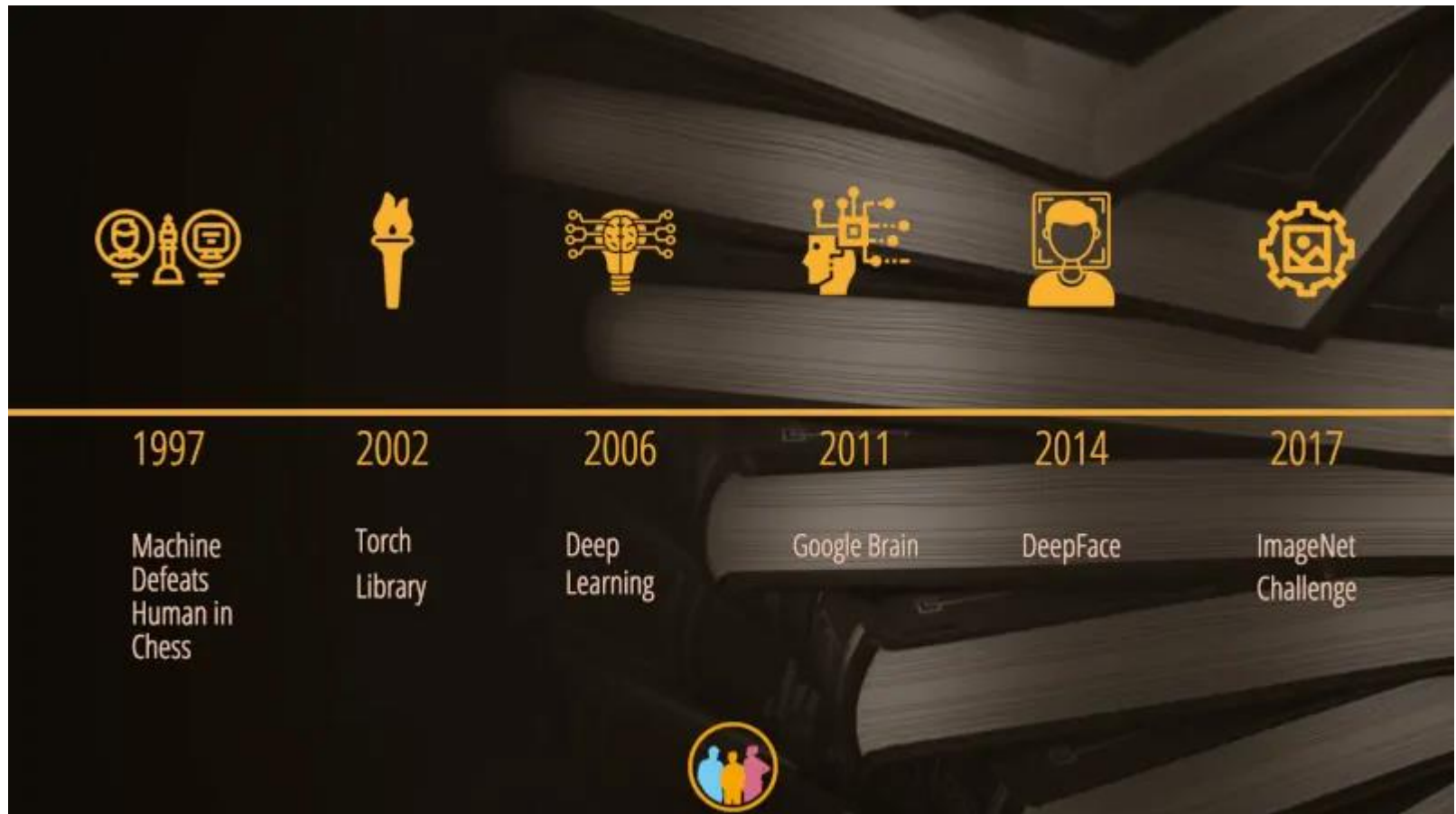
- Învățare automată (*învățarea mașinilor*)
- 1959, Arthur Samuel (IBM)



[Sursă](#)

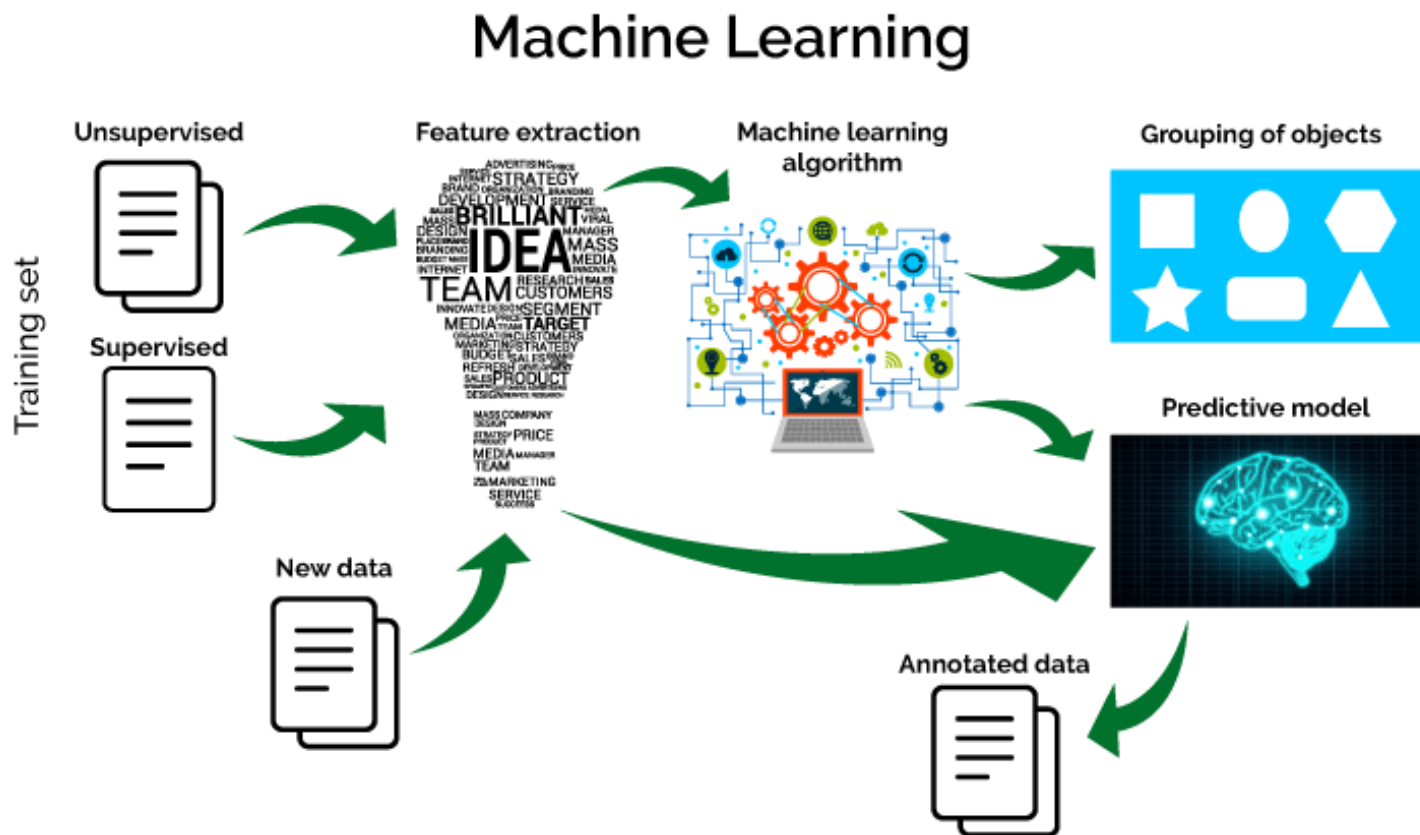
Machine Learning

- Moore's Law
- AI winter (acum suntem în AI Spring)



[Sursă](#)

Machine Learning

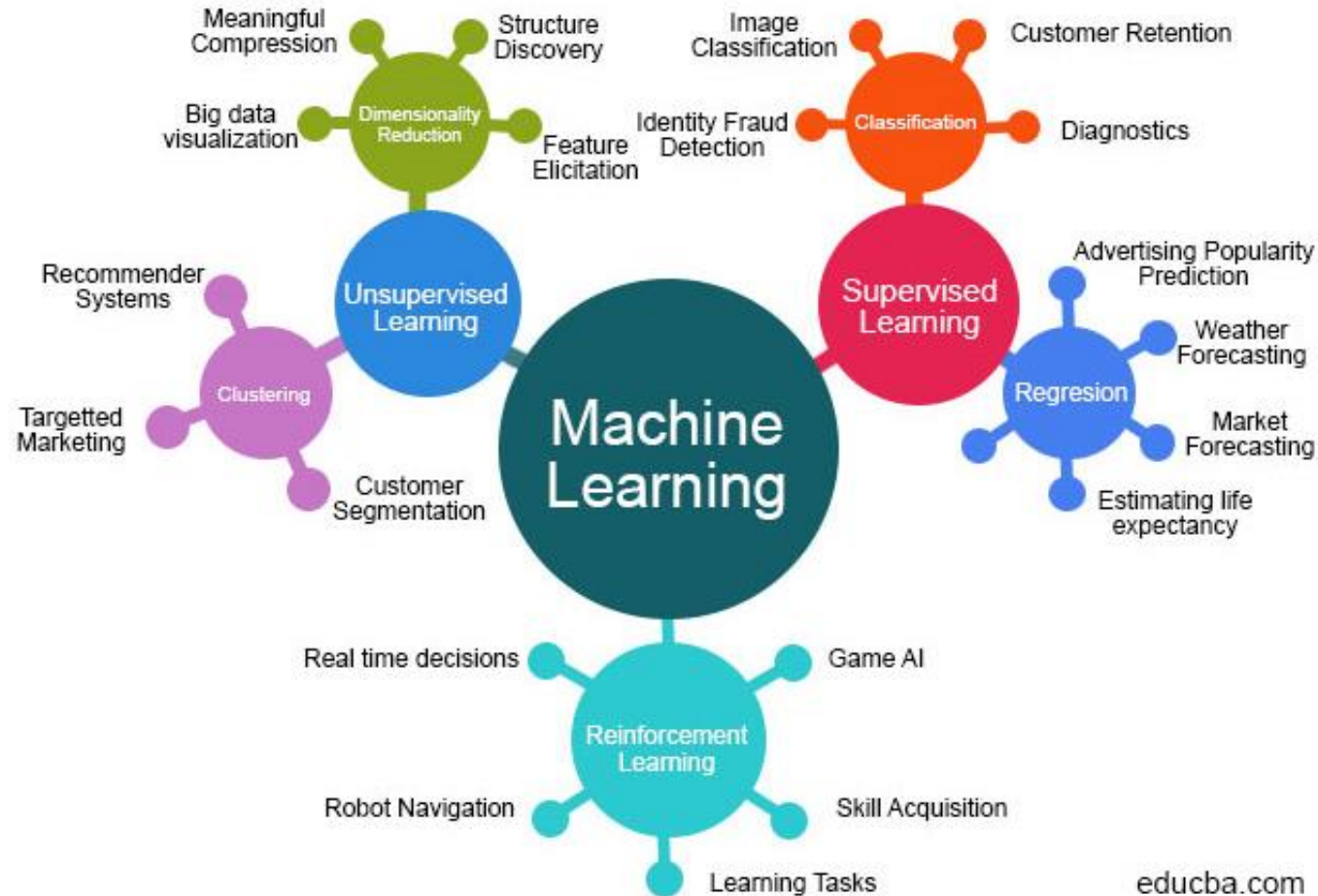


[Sursă](#)

Cine face extragerea trăsăturilor?

Machine Learning

➤ Aplicații



Clasificare

– date categorice

Regresie

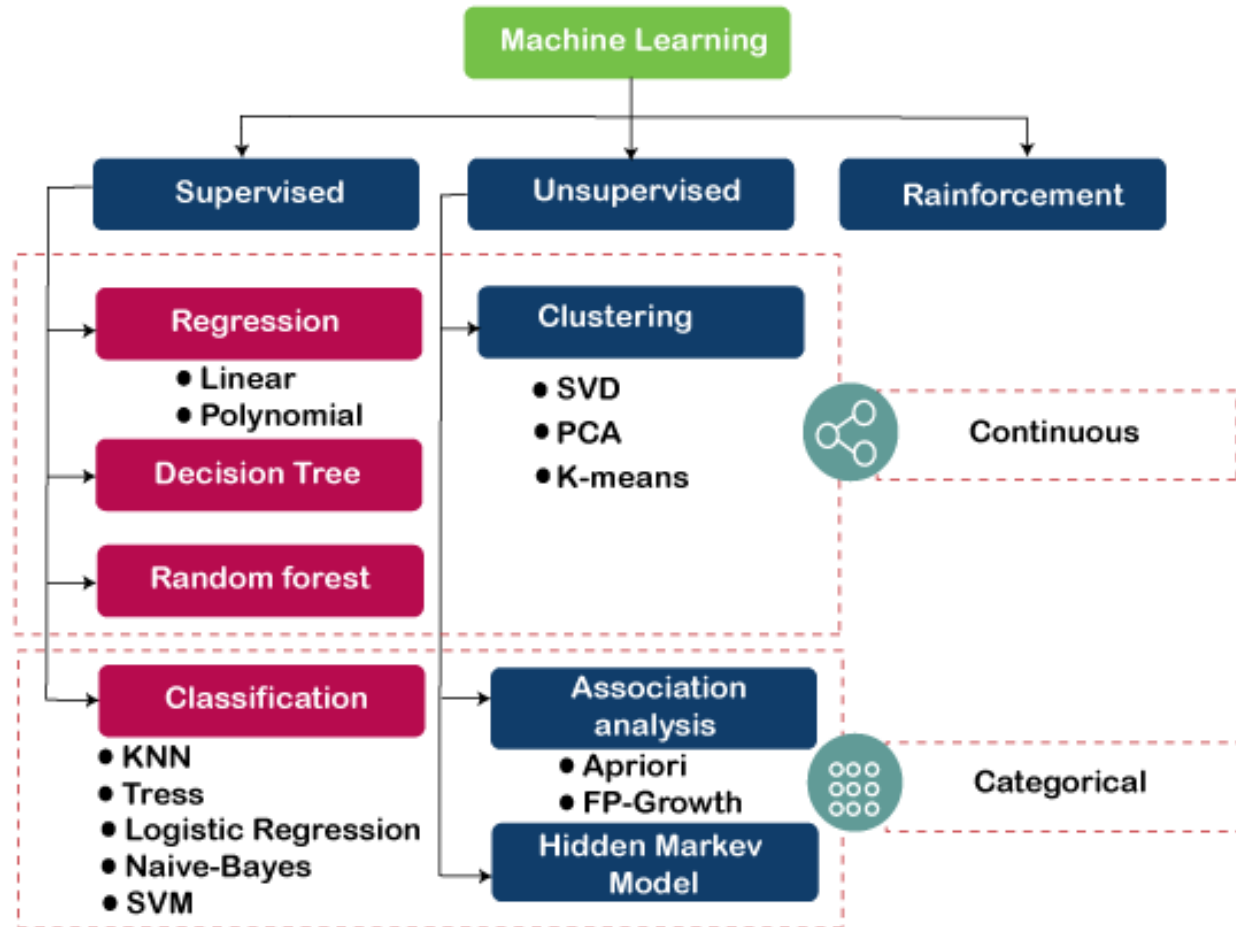
– date continue

educba.com

[Sursă](#)

Machine Learning

➤ Algoritmi



KNN – K Nearest Neighbour

SVM – Support Vector Machine

SVD – Singular Value Decomposition

PCA – Principal Component Analysis

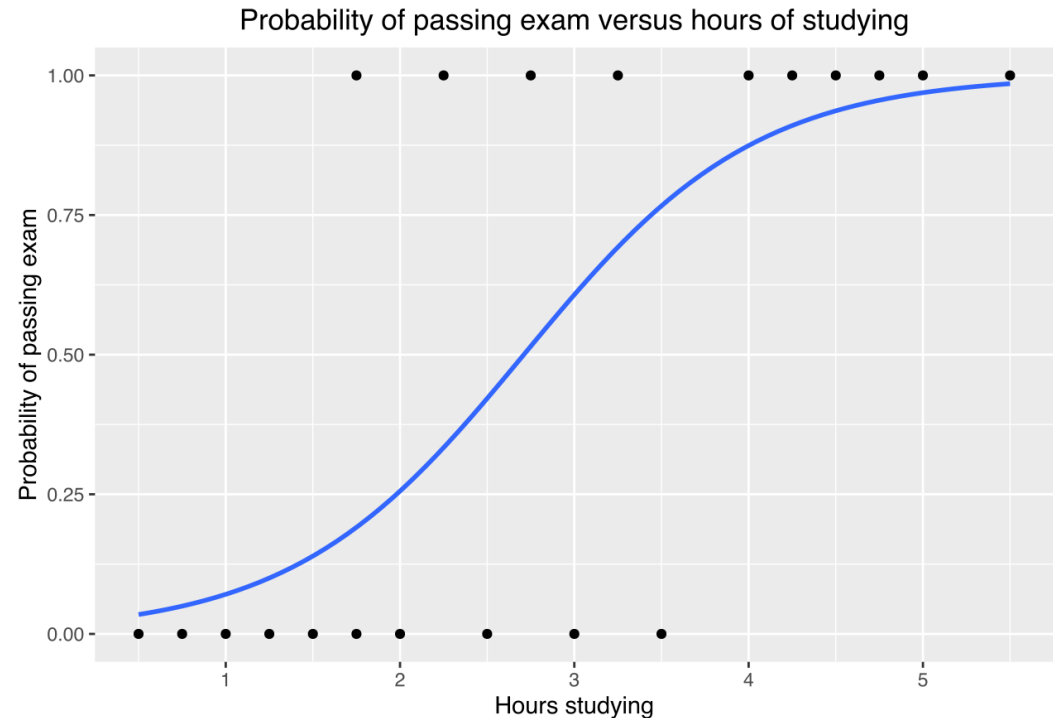
[Sursă](#)

Machine Learning

➤ Algoritmi – Logistic Regression

- pentru clasificare binară

$$f(x) = \frac{1}{1 + e^{-x}}$$

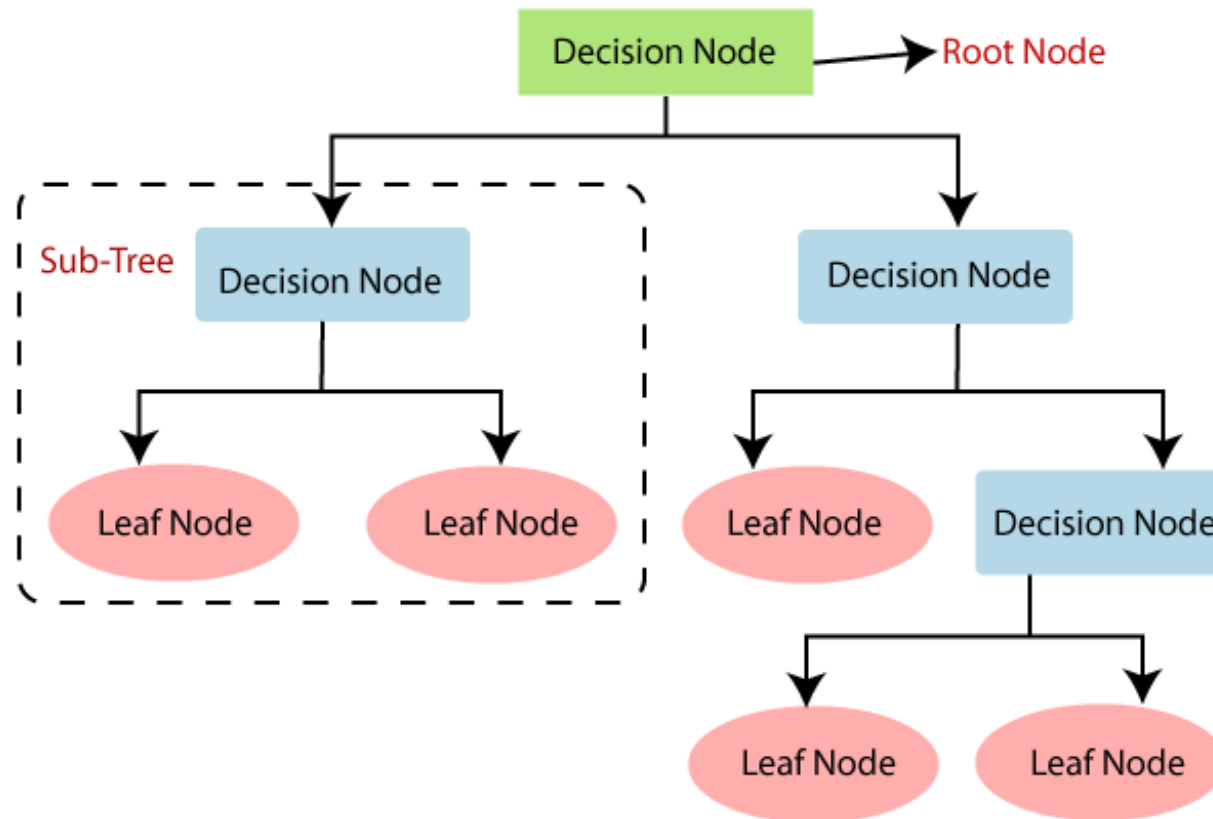


Hours (x_k)	0.50	0.75	1.00	1.25	1.50	1.75	1.75	2.00	2.25	2.50	2.75	3.00	3.25	3.50	4.00	4.25	4.50	4.75	5.00	5.50	
Pass (y_k)	0	0	0	0	0	0	1	0	1	0	1	0	1	0	1	1	1	1	1	1	1

[Sursă](#)

Machine Learning

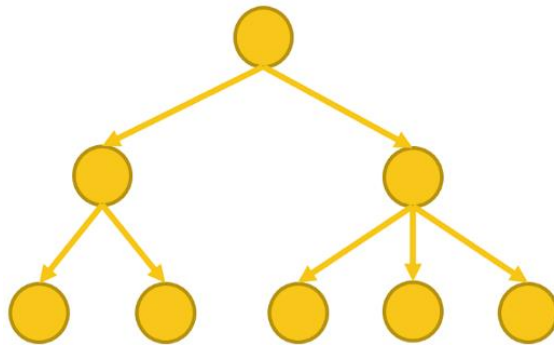
➤ Algoritmi – Decision Tree

[Sursă](#)

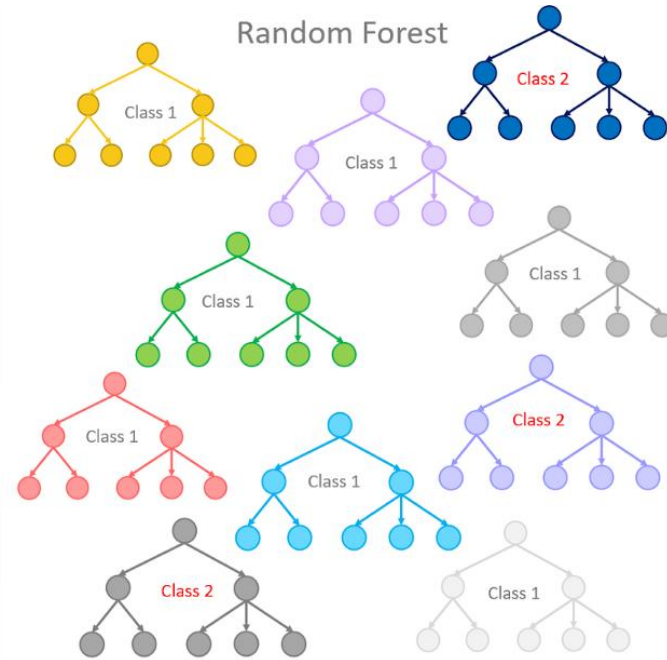
Machine Learning

➤ Algoritmi – Random Forest

Single Decision Tree



Random Forest



- N arbori de decizie sunt antrenati ușor diferit, iar ieșirea fiecăruia este luată în calcul pentru decizia finală
- predicția rezultată este mai precisă decât la arbori unici.

[Sursă](#)

Machine Learning

➤ Algoritmi – Naive Bayes

- clasificare probabilistică
- pentru fiecare obiect, se consideră că trăsăturile (*features*) contribuie independent la alocarea unui obiect la o clasă

Ex. Obiect = măr

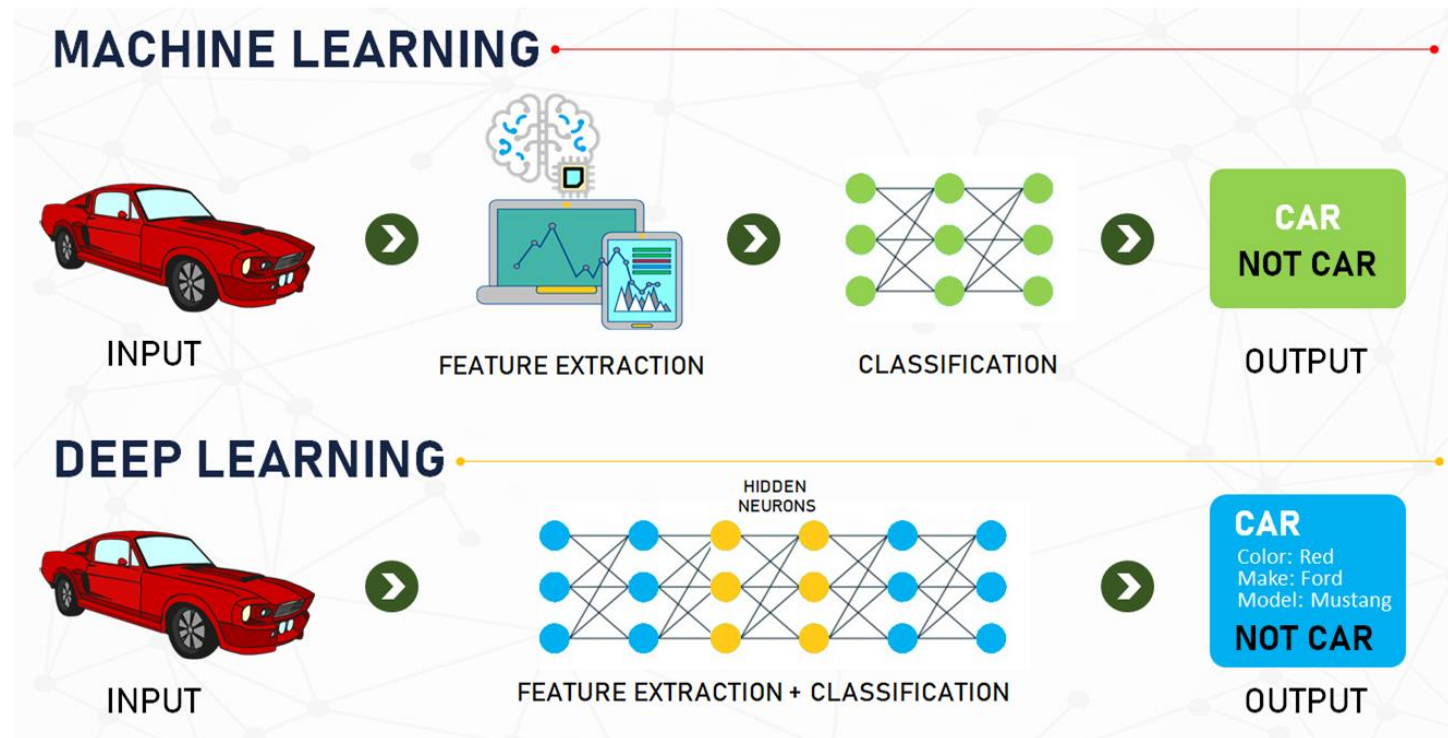
Trăsături: roșu, rotund, diametru 8 cm

Algoritmul Naive Bayes presupune că cele 3 trăsături nu sunt corelate.

[Sursă](#)

Machine Learning

➤ Limitări



- intervenție umană la etichetarea datelor

[Sursă](#)

Machine Learning

➤ Limitări



[Sursă](#)

- dimensiunea setului de date și calitatea acestora

Ex. Sistem de recrutare automată la Amazon (oprit în 2017)

- recomanda cu precădere candidați de sex masculin

- Clasificare – definire, tipuri ✓
- Etape ale clasificării ✓
- Machine Learning – definire, algoritmi, utilizare, limitări ✓

În episodul următor: **Inteligența Artificială – fundamente, terminologie, paradigme**