

Arbori de decizie: Privire de ansamblu

A. Noțiuni preliminare

- partiție a unei mulțimi;
- entropie (definiție: TM, pag. 57): pr. 2a, pr. 34a;
- entropie condițională specifică: pr. 14a;
- entropie condițională medie: pr. 2cd, ...;
- câștig de informație (definiție: TM, pag. 58): pr. 2cd, pr. 34c, pr. 36a;
- **arbori de decizie**, văzuți **ca structură de date**: pr. 1, pr. 6b, pr. 28
și, respectiv, ca program în logica propozițiilor: pr. 2e, pr. 33bc;
expresivitatea arborilor de decizie cu privire la funcții boolene: pr. 29;
- **spațiu de versiuni** pentru un concept (de învățat): pr. 1, pr. 28, pr. 32;
- **zone de decizie** și **granițe de separare/decizie** pentru arbori de decizie cu variabile continue: pr. 9.

B. Algoritmul ID3 și variante

B1. Algoritmul ID3

- pseudo-cod: TM, pag. 56;
- **bias-ul inductiv**: TM, pag. 63-64;
- exemple simple de aplicare: pr. 2, pr. 32, pr. 34;
în prezența „zgomotelor”: pr. 36;
- ID3 ca algoritm *per se*:
 - de tip divide-et-impera, recursiv,
greedy \Rightarrow nu garantează obținerea soluției optime:
pr. 4, pr. 21a, pr. 33 (vs. pr. 32b);
 - **1-step look-ahead**
 - complexitate de timp (vezi Weka book, 2011, pag. 199):
la antrenare, în anumite condiții: $\mathcal{O}(dm \log m)$; la testare $\mathcal{O}(d)$,
unde d este numărul de atribute, m este numărul de exemple
- ID3 ca algoritm ML:
 - **consistent** cu datele de antrenament (în absența „zgomotelor”): pr. 4
 - algoritm de învățare de tip **“eager”**
 - **analiza erorilor**:
la antrenare: pr. 5, pr. 6a, pr. 9a, pr. 38;
la validare
la n -fold cross-validare
la cross-validare leave-one-out (CVLOO): pr. 9b, pr. 40bc;
 - **overfitting**: pr. 9, pr. 21bc, pr. 40, pr. 57b.

B2. Extensii / variante ale algoritmului ID3

- **attribute cu valori continue:** pr. 9, pr. 14c, pr. 40, pr. 41; ch. 4, pr. 11b;
3-way splitting (sau: n -way splitting): pr. 12, pr. 42
alte variante de partiționare: pr. 43
- attribute discrete cu multe valori: pr. 13
- attribute cu valori nespecificate pentru unele instanțe;
- attribute cu diferite costuri asociate: pr. 14
- reducerea caracterului “eager” al învățării: pr. 16
- reducerea caracterului “greedy” al învățării:
IG cu “2-step look-ahead”: pr. 17, pr. 18
variante de tip “look-ahead” specifice atributelor continue: pr. 44
- folosirea altor măsuri de impuritate în locul câștigului de informație:
Gini Impurity, Misclassification Impurity: pr. 15
- reducerea **overfitting**-ului:
reduced-error pruning (folosind un set de date de validare):
TM, pag. 69-71; A. Cornuéjols, L. Miclet, 2nd ed., pag. 418-421;
rule post-pruning: TM, pag. 71-72;
top-down vs. bottom-up pruning: pr. 19, pr. 45;
pruning folosind testul statistic χ^2 : pr. 20, pr. 46

C. Proprietăți numerice / calitative ale arborelui ID3

- eroarea la antrenare produsă de algoritmul ID3 pe orice set de date consistente este 0;
- arborele ID3 nu este neapărat optimal (ca nr. de noduri/niveluri): pr. 4, pr. 33
- influența atributelor identice și, respectiv, a instanțelor multiple asupra arborelui ID3: pr. 7;
- o margine superioară pentru numărul de noduri frunză din arborele ID3, în funcție de numărul de exemple și de numărul de attribute: pr. 8;
- o margine superioară pentru eroarea la antrenare, în funcție de numărul de valori ale variabilei de ieșire): pr. 6b;
- o margine superioară pentru adâncimea arborelui ID3 când attributele de intrare sunt continue, iar datele de antrenament sunt (ne)separabile liniar: pr. 10;
- o aproximare simplă a numărului de instanțe greșit clasificate din totalul de M instanțe care au fost asignate la un nod frunză, cu ajutorul entropiei (H) nodului respectiv: pr. 37.

D. Alte metode de învățare automată bazate pe arbori

- metode de învățare automată de tip ansamblist bazate pe arbori de decizie:
boosting (AdaBoost): pr. 22 (convergența erorii de antrenare), pr. 23, 48, 49, 50, 51, 53 (aplicare), pr. 24 (AdaBoost ca algoritm de *optimizare secvențială* în raport cu funcția de cost / „pierdere” negativ-exponențială), pr. 25 (marginea

- de votare), pr. 26 (o condiție suficientă pentru γ -slab învățabilitate, bazată pe marginea de votare), pr. 52 (selectarea trăsăturilor folosind AdaBoost; aplicare la clasificarea de documente)), pr. 54 (o proprietate interesantă: orice mulțime de instanțe distincte [și etichetate] din \mathbb{R} este γ -slab învățabilă cu ajutorul compașilor de decizie), pr. 55 (o variantă generalizată a algoritmului AdaBoost), pr. 22 și 23 (adevărat / fals); Bagging, Random Forests;
- arbori de regresie (CART).

Clasificare bayesiană: Privire de ansamblu

A. Noțiuni preliminare (vezi [și] cap. Probabilități și statistică)

- probabilități și probabilități condiționate;
- formula lui Bayes: pr. 4b;
cap. Probabilități și statistică, pr. 6-7, 43-44;
- independența [condițională a] evenimentelor aleatoare:
cap. Probabilități și statistică, pr. 5, 40, 41;
- variabile aleatoare:
estimare de parametri în sensul verosimilității maxime (MLE): pr. 4a;
(vezi și cap. Estimarea probabilităților, pr. 5, 20)
- distribuții probabiliste corelate, marginale și condiționale: pr. 6, 8, 10, 27.
(vezi și cap. Probabilități și statistică, pr. 14, 15)
- independența [condițională a] variabilelor aleatoare: pr. 8, 10, 24, 27-33;
(vezi și cap. Probabilități și statistică, pr. 16, 53, 19, 20, 62)
- ipoteze MAP vs. ipoteze ML:
formulare [ca soluții la] probleme de optimizare: TM, pag. 156-157 (vezi pr. 13);
exemplificare: pr. 2, 3, 21, 1, 32;
exemplificare în cazul arborilor de decizie: pr. 12;
- regresia logistică: vezi draftul capitolului suplimentar pentru cartea ML a lui T. Mitchell, *Generative and discriminative classifiers: Naive Bayes and logistic regression* (în special secțiunea 3);
- regresia logistică, chestiuni introductive: pr. 10 de la capitolul *Estimarea probabilităților*.

B. Algoritmi

- Algoritmul Bayes Naiv și algoritmul Bayes Corelat:³²⁷
formulare ca probleme de optimizare: TM, pag. 167;
pseudo-cod: vezi slide-uri;
aplicare: pr. 4, 6, 7, 22, 23, 24;
- aplicarea/adaptarea algoritmului Bayes Naiv pentru clasificare de texte:³²⁸
pr. 5, 25
folosirea regulii “add-one” [a lui Laplace] pentru „netezirea” parametrilor:
pr. 5, 26;
- calculul ratei medii a erorilor pentru algoritmii Bayes Naiv și Bayes Corelat:
pr. 8, 9, 27, 28, 29, 33;
- zone de decizie și granițe de decizie (separatori):
evidențierea grafică a „erorilor” clasificatorului Bayes Naiv în raport cu clasificatorul Bayes Corelat: pr. 10.

³²⁷La punctele B și C considerăm (implicit) că toate variabilele de intrare sunt de tip Bernoulli sau, mai general, de tip categorial. La punctul D vom considera și variabile de intrare de tip continuu, în genere de tip gaussian. Variabila de ieșire se consideră întotdeauna de tip Bernoulli/categorial.

³²⁸Atenție: Noi am folosit aici versiunea de bază a algoritmului Bayes Naiv; varianta “bag of words” (vezi cartea Machine Learning a lui Tom Mitchell, pag. 183) diferă ușor de aceasta.

C. Proprietăți ale algoritmilor Bayes Naiv și Bayes Corelat

- dacă proprietatea de independență condițională a atributelor de intrare în raport cu variabila de ieșire se verifică, atunci rezultatele produse de către cei doi algoritmi (Bayes Naiv și Bayes Corelat) în faza de testare coincid;
- numărul de parametri necesari de estimat din date: liniar pentru Bayes Naiv ($2d + 1$) și exponențial pentru Bayes Corelat ($2^{d+1} - 1$): pr. 7e, 23ab, 28;
- complexitatea algoritmului Bayes Naiv:
 - complexitatea de spațiu: $\mathcal{O}(dn)$
 - complexitatea de timp:
 - la antrenare: $\mathcal{O}(dn)$
 - la testare: $\mathcal{O}(d')$,
 - unde n este numărul de exemple, iar d este numărul de atribute de intrare [LC: d' este numărul de atribute de intrare din instanța de test]
- complexitatea de eșantionare: de ordin logaritm pentru Bayes Naiv și de ordin exponențial pentru Bayes Corelat: pr. 11.
- (P0) echivalența regulei de decizie a algoritmului Bayes Naiv (când toate variabilele de intrare sunt de tip Bernoulli) regulei de decizie a *regresiei logistice* și, în consecință, liniaritatea granițelor de decizie.

D. Algoritmii Bayes Naiv și Bayes Corelat cu variabile de intrare de tip gaussian

Proprietăți:

- (P1) presupunem că variabila de ieșire este booleană, i.e. ia valorile 0 sau 1; dacă pentru orice atribut de intrare, variabilele condiționale $X_i|Y = 0$ și $X_i|Y = 1$ au distribuții gaussiene de varianțe egale ($\sigma_{i0} = \sigma_{i1}$), atunci regula de decizie GNB (Gaussian Naive Bayes) este echivalentă (ca formă) cu cea a regresiei logistice, deci separarea realizată de către algoritmul GNB este de formă liniară (vezi pr. 34a, 15);

Aplicare: G[N]B: 16, 38, G[J]B: 37, GNB vs G[J]B: 18.

E. Alte chestiuni

- comparații între algoritmul Bayes Naiv și alți algoritmi de clasificare automată: pr. 31, 33;
- calculul [intervalului] „erorii reale” pornind de la „eroarea de eșantionare” (eroarea de test) pentru un clasificator oarecare: cap. Probabilități și statistică, pr. 27.

Învățare bazată pe memorare: Privire de ansamblu

A. Noțiuni preliminare

- măsuri de distanță, măsuri de similaritate: pr. 2
- normă într-un spațiu vectorial; [măsura de] distanță indusă de către o normă: pr. 7

B. Algoritmul k -NN și variante

B1. Algoritmul k -NN

- pseudo-cod: TM, pag. 232
- bias-ul inductiv: ... „Cine se aseamănă se adună“ (sau: „Spune-mi cu cine te împrietenești, ca să-ți spun cine ești“)
- exemple (simple) de aplicare: pr. 1, pr. 2
- complexitate de spațiu: $\mathcal{O}(dn)$
complexitate de timp:
 - la antrenare: $\mathcal{O}(dn)$
 - la testare: $\mathcal{O}(dn \log n)$

[LC: $\mathcal{O}(dnk \log k)$ pt. $k > 1$ (worst case) și $\mathcal{O}(dn)$ pt. $k = 1$],

unde d este numărul de atribute, iar n este numărul de exemple
- arbori kd (engl., kd -trees): *Statistical Pattern Recognition*, pag. 163-173
- k -NN ca algoritm ML “lazy” (vs. “eager”):
suprafețe de decizie și granițe de decizie:
diagrame Voronoi pentru 1-NN: pr. 4, pr. 11, pr. 17, pr. 18, pr. 19;
Proprietate: suprafețele de decizie și granițele de decizie depind de măsura de distanță folosită: pr. 7
- analiza erorilor:
 - 1-NN pe date consistente: eroarea la antrenare: 0
 - variația numărului de erori (la testare și respectiv testare) în funcție de valorile lui k : pr. 20, pr. 21ab
 k -NN ca metodă ne-parametrică; alegerea lui k : CV: pr. 21c
 - CVLOO: pr. 3, pr. 12, pr. 15, pr. 19
 - sensibilitatea / robustețea la „zgomote”: pr. 5
 - eroarea asimptotică: pr. 10, pr. 23
- efectul trăsăturilor redundante sau irelevante
- „blestemul marilor dimensiuni” (engl., the curse of dimensionality): pr. 9.

B2. Variante ale algoritmului k -NN

- k -NN folosind alte măsuri de distanță (decât dist. euclidiană): pr. 7
- k -NN cu ponderarea distanțelor (engl., distance-weighted k -NN):
TM, pag. 236-238 (form. 8.2, 8.3, 8.4)
- algoritmul lui Shepard: pr. 8

C. Comparații cu alți algoritmi

- ID3: pr.11, pr. 13ab
- SVM: pr.12, pr. 13c

D. Alte metode de tip IBL

- regresie local-ponderată: TM, pag. 236-238
- rețele RBF: TM, pag. 238-240
- raționare bazată pe cazuri (engl., case-based reasoning): TM, pag. 240-244

Clusterizare: Privire de ansamblu

A. Noțiuni de bază

- instanță neetichetată vs. instanță etichetată (exemplu de antrenament)
- învățare nesupervizată (clusterizare) vs. învățare supervizată (clasificare)
- cluster / grup / grupare / bin (engl.) vs. clasă
- [funcție/măsură de] *distanță* definită pe $\mathbb{R}^d \times \mathbb{R}^d$
- tipuri de clusterizare: ierarhică vs. neierarhică
- tipuri de ierarhii: ierarhii (arbori de clusterizare, dendrograme) obișnuite vs. ierarhii plate (engl., flat hierarchies);
exemple: pr. 1a și respectiv pr. 1b, pr. 6a
- tipuri de apartenență a unei instanțe la un cluster: hard vs. soft (ultima numai pt. clusterizare ne-ierarhică)

B. Clusterizare ierarhică

B1. Noțiuni specifice

- [funcție de] *similaritate* între clustere, definită pe baza [extinderii] noțiunii de distanță la $\mathcal{P}(X) \times \mathcal{P}(X)$, unde $X \subset \mathbb{R}^d$ este mulțimea de instanțe, iar $\mathcal{P}(X)$ este mulțimea părților lui X ;

tipuri de [funcții de] similaritate:

“single-linkage”:³²⁹ $d(A, B) = \min\{d(x, y) | x \in A, y \in B\}$

“complete-linkage”:³³⁰ $d(A, B) = \max\{d(x, y) | x \in A, y \in B\}$

“average-linkage”: $d(A, B) = \frac{1}{|A| |B|} \sum_{x \in A, y \in B} d(x, y)$.

În general, putem considera $\text{sim}(A, B) = 1/(1 + d(A, B))$ sau chiar $\text{sim}(A, B) = 1/d(A, B)$ când ne referim doar la clustere non-singleton;

proprietate/restricție: $\text{sim}(A \cup B, C) \leq \min\{\text{sim}(A, C), \text{sim}(B, C)\}$ pentru orice clustere A, B selectate de algoritmul de clusterizare ierarhică la un pas oarecare și orice alt cluster C ;

- [funcție de] *coeziune* [internă] a unui cluster (sau: între elementele / instanțele dintr-un cluster);

exemplu (pentru clustere non-singleton):

$$\text{coh}(A) = \left(\frac{1}{C_{|A|}^2} \sum_{x, y \in A} d(x, y) \right)^{-1} = \frac{C_{|A|}^2}{\sum_{x, y \in A} d(x, y)}$$

³²⁹Sau: nearest-neighbour.

³³⁰Sau: furthest-neighbour.

B2. Algoritmi de clusterizare ierarhică

- tipuri de algoritmi de clusterizare ierarhică:
 - bottom-up (*clusterizare aglomerativă*) vs. top-down (*clusterizare divizivă*);
 - pseudo-cod: Manning & Schütze, *Foundations of Statistical Natural Language Processing*, 2002, pag. 502;
 - analiza (ca algoritmi *per se*): ambii algoritmi sunt iterativi și “greedy”; rezultatele (ierarhiile) obținute nu sunt determinate neapărat în mod unic (vezi pr. 3b);
 - exemple de aplicare: pr. 1-5, 25-29 (pentru bottom-up), respectiv pr. 6 (pentru top-down);
 - implementări: pr. 54, 31, 55.
- Proprietăți:
 - clusterizarea folosind similaritate de tip “single-linkage” are tendința să creeze clustere alungite; invers, folosind similaritate “complete-linkage” sau “average-linkage”, se formează clustere de formă mai degrabă sferică (vezi pr. 5 și 28);
 - numărul maxim de niveluri dintr-o dendrogramă (văzută ca arbore în sensul teoriei grafurilor) este $n - 1$, unde n este numărul de instanțe de clusterizat (vezi pr. 4a); numărul minim de niveluri: $\lceil \log_2 n \rceil$ (vezi pr. 4b).
 - există o anumită corespondență între clusterizare ierarhică cu similaritate de tip
 - “single-linkage” și aflarea *arborelui [de acoperire] de cost minim* dintr-un graf; vezi pr. 6;
 - “complete-linkage” și aflarea unei *clici* (subgraf maximal complet) dintr-un graf; (vezi Manning & Schütze, *op. cit.*, pag. 506-507)
 - algoritmul de clusterizare aglomerativă la al cărui pseudo-cod am făcut referire mai sus are complexitate $\mathcal{O}(n^3)$ (vezi pr. 25); atunci când se folosește single-linkage sau complete-linkage, există însă versiuni/algoritmi de complexitate $\mathcal{O}(n^2)$: SLINK (1973) și respectiv CLINK (1976);
 - la clusterizare ierarhică aglomerativă cu similaritate “average-linkage”:

dacă se folosește ca măsură de similaritate între 2 instanțe cosinusul unghiului dintre vectorii care instanțele și se „normalizează” acești vectori (i.e., se lucrează cu 2 vectori coliniari cu ei, dar de normă egală cu 1), atunci calculul coeziunii [interne a] unui cluster nou format, precum și calculul „distanței” dintre două clustere se pot face în timp constant 32.
- Alte tipuri de măsuri de similaritate:

metrica lui Ward: pr. 30

C. Clusterizare neierarhică...

C1. ...folosind asignare “hard” a instanțelor la clustere

Noțiuni specifice

- centroid (centru de greutate) al unui cluster,
 K -partiție, K -configurație [inițială] a centroizilor (pr. 11);
- o funcție de evaluare a „calității” clusterelor (sau: funcție de „coeziune” / „distorsiune” / „eroare” totală):
 „suma celor mai mici pătrate”: $J_K(C, \mu) = \sum \|x_i - \mu_{C(x_i)}\|^2$, unde C este K -partiție, μ este K -configurație de centroizi, iar $\mu_{C(x_i)}$ este centroidul cel mai apropiat de x_i (pr. 12).

Algoritmul K -means

- pseudo-cod (o versiune [mai] generală): Manning & Schütze, *op. cit.*, pag. 516; alternativ, vezi enunțul pr. 12 (sau, echivalent, folosind variabile-indicator: pr. 37);
 exemple de aplicare: pr. 7-11, 15a, 19a, 20a, 33, 34;
- exemple de *euristici pentru inițializarea centroizilor*:
 inițializare arbitrară/random în \mathbb{R}^d sau în $\{x_1, x_2, \dots, x_n\} \in \mathbb{R}^d$ (setul de date de clusterizat);
 aplicare în prealabil a unui algoritm de clusterizare ierarhică;
 K -means++ (David Arthur, Sergei Vassilvitskii, 2007).
- *exemple de criterii de oprire*:
 după efectuarea unui număr maxim de iterații (fixat inițial);
 când componența clusterelor nu se mai modifică de la o iterație la alta;
 când pozițiile centroizilor nu se mai modifică de la o iterație la alta;
 când descreșterea valorii criteriului J_K de la o iterație la alta nu mai este strictă sau nu mai este peste un anumit prag ε fixat în prealabil.
- ca algoritm *per se*:
algoritm iterativ: pleacă de la o soluție (K -partiție) aleasă eventual în mod arbitrar/aleatoriu și o „îmbunătățește” la fiecare iterație;
 soluția găsită este dependentă de inițializarea centroizilor (vezi pr. 10);
 mai mult, chiar la o aceeași inițializare, rezultatele pot diferi(!) dacă avem instanțe multiple/redundante, situate la egală distanță de 2 centroizi la o iterație oarecare (vezi pr. 12b);
 K -means poate fi văzut și ca *algoritm de optimizare* — vezi criteriul J_K de mai sus;
 explorează doar parțial spațiul de căutare [a minimului criteriului J_K];
 algoritmul K -means nu garantează atingerea optimului global (i.e., minimul criteriului J_K (vezi pr. 12b, 38b);
 strategia de căutare/optimizare folosită de K -means este de tipul *descreștere pe coordonate* (engl., coordinate descent), i.e. descreștere iterativă, mergând alternativ pe fiecare din cele două coordonate ale criteriului $J_K(C^t, \mu^t)$ (vezi pr. 12a).

- ca algoritm de învățare automată:
[urmat de] „generalizare”: o instanță nouă x se asociază clusterului având centroidul cel mai apropiat de x ;
„granițele” de separare dintre [perechile de] clustere produse de K -means sunt [doar] liniare, [cel puțin] atunci când se folosește distanța euclidiană (vezi pr. 12.b);
este însă posibil să se obțină separatori neliniari dacă se folosește o versiune „kernelizată” a algoritmului K -means (vezi pr. 39);
ca *euristică* pentru alegerea unei valori convenabile pentru K , vezi CMU, 2012 fall, E. Xing, A. Singh, HW3, pr. 1.
- implementare: pr. 56.

Proprietăți

- în legătură cu criteriul definit mai sus, $J_K : \mathcal{P}_K \times (\mathbb{R}^d)^K \leftarrow [0, +\infty)$, unde \mathcal{P}_K este mulțimea tuturor K -partițiilor peste mulțimea de instanțe, $X = \{x_1, x_2, \dots, x_n\} \subseteq \mathbb{R}^d$:
 - valoarea 0 este atinsă, și anume atunci când $K = n$, C este K -partiția de clustere singleton $C_i = \{x_i\}$, iar $\mu_i = x_i$, pentru $i = 1, \dots, n$ (pr. 36);
 - pentru $K > 0$ fixat, $|\mathcal{P}_K| = K^n$, deci este finit, și există $\underline{J}_K \stackrel{\text{not.}}{=} \min_C J_K(C, \mu_C)$; acest minimum (\underline{J}_K) se poate obține prin explorarea exhaustivă a spațiului \mathcal{P}_K , însă consumul de timp este prohibitiv în practică (vezi pr. 12a);
 - $\underline{J}_1 \geq \underline{J}_2 \geq \dots \geq \underline{J}_{n-1} \geq \underline{J}_n = 0$ (vezi pr. 13).
- în legătură cu J_K și algoritmul K -means:
 - $J_K(C^{t-1}, \mu^{t-1}) \geq J_K(C^t, \mu^t)$ la orice iterație ($t > 0$) a algoritmului K -means (vezi pr. 12a);
 - în consecință, dacă se impune restricția ca la fiecare iterație inegalitatea de mai sus să fie satisfăcută în varianta strictă ($J_K(C^{t-1}, \mu^{t-1}) > J_K(C^t, \mu^t)$), atunci algoritmul K -means termină într-un număr finit de pași;
 - în vreme ce minimizează *coeziunea intra-clustere*, i.e. o sumă ponderată a „sumelor celor mai mici pătrate” calculate pe clustere,

$$\sum_{k=1}^K \frac{\sum_{i=1}^n \gamma_{ik} \|x_i - \mu_k\|^2}{\sum_{i=1}^n \gamma_{ik}}$$

unde $\gamma_{ik} = 1$ dacă x_i aparține clusterului de centroid μ_k și $\gamma_{ik} = 0$ în caz contrar, algoritmul K -means maximizează (în mod aproximativ!) o măsură de distanță între clustere

$$\sum_{k=1}^K \left(\frac{\sum_{i=1}^n \gamma_{ik}}{n} \right) \|\mu_k - \bar{x}\|^2$$

unde \bar{x} este media instanțelor x_1, x_2, \dots, x_n (pr. 37).

- dacă $d = 1$, deci $x_1, x_2, \dots, x_n \in \mathbb{R}$,
 - orice K -partiție (C_1, \dots, C_K) pentru care se atinge \underline{J}_K este de forma unei colecții de „intervale”: $C_1 = \{x_1, \dots, x_{i_1}\}$, $C_2 = \{x_{i_1+1}, \dots, x_{i_2}\}$, ..., $C_{K-1} = \{x_{i_{K-1}+1}, \dots, x_n\}$, cu $i_1 < i_2 < \dots < i_{K-1} < i_K = n$;
 - există un algoritm de complexitate $\mathcal{O}(Kn^2)$ care calculează \underline{J}_K (vezi pr. 38).

- algoritmul K -means poate fi „kernelizat” (vezi pr. 39); în consecință, putem obține drept granițe de separare între clustere [și] suprafețe non-liniare (spre deosebire de versiunea ne-kernelizată a algoritmului K -means, unde granițele sunt doar liniare).

C2. ...folosind asignare “soft” a instanțelor la clustere

Noțiuni preliminare

- variabile aleatoare (discrete, resp. continue); media, varianța și co-varianța variabilelor aleatoare;
- vector de variabile aleatoare; matrice de covarianță pentru un astfel de vector; proprietăți: matricea de covarianță trebuie să fie în mod necesar simetrică și pozitiv definită (vezi pr. 24 de la capitolul de *Probabilități și statistică*);
- distribuție (funcție de densitate) de probabilitate (p.d.f.); parametri ai unei distribuții; distribuția gaussiană: cazurile uni- și multi-variant;
- mixtură de distribuții probabile: văzută ca o formă particulară de *combinație liniară* de distribuții de probabilitate $\pi_1\Psi_1 + \pi_2\Psi_2 + \dots + \pi_k\Psi_k$ (cu $\pi_i \geq 0$ și $\sum_{i=1}^k \pi_i = 1$), definită [și mai specific] scriind distribuția $P(X)$ ca o sumă ponderată de probabilități condiționate: $\sum_z P(X|Z)P(Z)$, unde X sunt variabilele „observabile”, iar variabila Z (eventual multiplă) poate fi „neobservabilă” / „latentă” / „ascunsă”; exemple: o mixtură de distribuții categoriale: pr. 56 de la capitolul *Probabilități și statistică*; o mixtură de distribuții gaussiene: pr. 16 de la capitolul *Clasificare bayesiană*.
- funcție de *verosimilitate* a unui set de date (D), în raport cu o distribuție probabilistă dată: $L(\theta) = P(D|\theta)$, unde prin θ se notează parametrii respectivei distribuții. Exemplificare: pr. 16c; pr. 1 de la cap. *Estimarea probabilităților*
- MLE (Maximum Likelihood Estimation): estimarea [valorilor] parametrilor unei distribuții probabile în sensul maximizării verosimilității datelor disponibile. Exemplificare: cap. *Estimarea probabilităților*, pr. 1-7, 11-21.
- Observație: Algoritmul EM este o metodă de estimare a parametrilor unei mixturi de distribuții probabile. *Alternativ*, pentru același obiectiv pot fi folosite alte metode, de exemplu *metoda gradientului ascendent*.

Algoritmul EM pentru clusterizare

prin estimarea parametrilor unui

model de mixturi de distribuții gaussiene (GMM)

- pseudo-cod: cazul uni-dimensional, varianta când doar parametrul μ este lăsat liber: *Machine Learning*, Tom Mitchell, 1997, pag. 193; aplicare: pr. 16; cazul uni-dimensional, varianta când toți parametrii (π , μ și σ) sunt lăsați liberi: pr. 17; alte variante: pr. 43, 44; cazul multi-dimensional, varianta când toți parametrii (π , μ și Σ) sunt lăsați liberi: pr. 23;

- schema algoritmică EM: *ML book*, pag. 195;
- ca algoritm *per se*:
 - *algoritm iterativ*: pleacă de la o soluție (instanțiere pentru parametri) aleasă eventual în mod arbitrar/aleatoriu și o „îmbunătățește” la fiecare iterație. Soluția găsită este dependentă de valorile inițiale ale parametrilor;
 - *algoritm de optimizare*: la fiecare iterație t se calculează o funcție „auxiliară” $Q_t(\theta|\theta^{(t)})$, care reprezintă media funcției de log-verosimilitate a datelor „complete” (cele „observabile” plus cele „neobservabile”), unde $\theta^{(0)}$, constând din valorile inițiale ale parametrilor mixturii (θ) se alege în mod arbitrar, iar apoi $\theta^{(t+1)} = \operatorname{argmax}_{\theta} Q_t(\theta|\theta^{(t)})$; media reprezentată de Q_t se calculează în funcție de distribuțiile condiționale ale variabilelor „neobservabile” în raport cu datele observabile și cu $\theta^{(t)}$; se poate demonstra că funcția Q_t constituie o margine inferioară pentru funcția de log-verosimilitate a variabilelor „observabile”, $\log P(X|\theta)$ (vezi pr. 1 de la capitolul *Algoritmul EM*); teorema de corectitudine / convergență (vezi problemele 1 și în special 2 de la capitolul *Algoritmul EM*) pe de o parte garantează faptul că la fiecare iterație a algoritmului EM, log-verosimilitatea datelor „observabile”, $\log P(X|\theta^{(t)})$ nu descrește (ci fie crește, fie rămâne neschimbată), dar pe de altă parte nu garantează găsirea optimului global al funcției de log-verosimilitate a datelor „observabile”, $\log P(X|\theta)$, ci eventual a unui optim local; metoda de căutare a optimului / maximumului funcției $\log P(X|\theta)$ este „coordinate ascent” (căutare pe coordonate, în mod alternant);
- ca algoritm de *învățare statistică*: algoritmul EM poate fi văzut ca o metodă de estimare a parametrilor (engl., parameter fitting);
- ca algoritm de *învățare automată*: algoritmul EM este o metodă de identificare/învățare de ipoteze ML (Maximum Likelihood); vezi capitolul/secțiunea 6.4 din cartea *Machine Learning*; învățare în prezența unor variabile aleatoare ne-observabile(!); [urmată eventual de] „generalizare”: o instanță nouă x se asociază clusterului (i.e., distribuției) j pentru care se atinge $\max_{j'} P(X = x|h_{j'})P(h_{j'})$; spre deosebire de cazul algoritmului K -means, suprafețele / granițele de separare create de algoritmul EM/GMM nu sunt în mod neapărat liniare (vezi de exemplu o situație întâlnită la rezolvarea pr. 15.c, pag. 358, sau pr. 51.c).
- comparativ cu algoritmul K -means: algoritmul EM/GMM este în general mai lent — mișcarea centrozilor poate explora într-o manieră mai fină spațiul (vezi de exemplu pr. 19) —, iar din acest motiv poate să obțină uneori rezultate mai bune / convenabile (vezi spre exemplu pr. 20), și este mai robust la influența outlier-elor; apare un fenomen de „atracție” reciprocă a mediilor gaussianelor (aceste medii fiind echivalentul centrozilor din algoritmul K -means), datorită faptului că fiecare instanță aparține (cu o anumită probabilitate) la fiecare cluster (vezi spre exemplu pr. 15.b).
- *schema algoritmică EM* (vezi Tom Mitchell, *Machine Learning book*, 1997, pag. 195) are diverse variante/aplicații:
 - calculul parametrilor pentru mixturi de diverse distribuții [nu doar gaussiene]: vezi capitolul *Algoritmul EM*;

- calculul parametrilor pentru gramatici probabiliste independente de context (engl., probabilistic context-free grammars, PCFG);
- calculul parametrilor modelelor Markov ascunse (engl., hidden Markov models, HMM);
- calculul parametrilor rețelelor bayesiene (engl., Bayes nets);
- calculul parametrilor rețelelor de funcții cu baza radială (engl., radial basis functions, RBF) o familie de rețele neuronale artificiale; etc.

Proprietăți

- Pentru distribuții gaussiene multi-variate:
 - dacă matricea de covarianță Σ este diagonală, atunci distribuția gaussiană respectivă este echivalentă cu un set/vector de variabile gaussiene uni-variate independente (vezi pr. 25 de la capitolul *Probabilități și statistică*);
 - dacă matricea Σ este de forma $\sigma^2 I$, unde I este matricea identitate, datele generate de respectiva distribuție tind să se grupeze în sfere;
 - dacă matricea Σ este diagonală (fără nicio altă restricție), datele generate se grupează în elipse (sau: corpuri elipsoidale) având axele de simetrie paralele cu axele sistemului de coordonate;
 - în cazul cel mai general (deci când matricea Σ nu este neapărat diagonală), datele generate de acest tip de distribuție se grupează în elipse (corpuri elipsoidale) cu axele de simetrie [desigur, perpendiculare, dar altfel] nerestricționate.
- Pentru schema algoritmică EM:
vezi cele menționate mai sus în legătură cu algoritmul EM văzut ca *algoritm de optimizare*.
- Legătura dintre algoritmul K -means și algoritmul EM/GMM (cazul multi-variat):
atunci când $\Sigma = \sigma^2 I$, iar $\sigma^2 \rightarrow 0$ (și sunt satisfăcute încă două restricții), algoritmul EM/GMM tinde să se comporte ca și algoritmul K -means (vezi pr. 47);
- O legătură interesantă între clasificatorul Bayes Naiv gaussian și algoritmul EM/GMM când matricele de covarianță sunt diagonale:
o variantă semi-supervizată a algoritmului EM/GMM. (pr. 53).

Algoritmul EM: Privire de ansamblu

Noțiuni preliminare

- estimarea parametrilor unei distribuții probabiliste în sensul verosimilității maxime (MLE) respectiv în sensul probabilității maxime a posteriori (MAP): vezi cap. *Probabilități și statistică*;
- tipuri/clase de distribuții probabiliste vezi cap. *Probabilități și statistică*;
- mixturi de distribuții probabiliste vezi cap. *Probabilități și statistică* și cap. *Clusterizare*;
- metoda “coordinate ascent” pentru rezolvarea problemelor de optimizare: pr. 1;
- metoda multiplicatorilor lui Lagrange pentru rezolvarea problemelor de optimizare cu restricții: pr. 5, 7 și 15.

Schema algoritmică EM

- pseudo-cod: *Machine Learning*, Tom Mitchell, 1997, pag. 194-195;
- fundamentare teoretică: pr. 1 și 2;
- chestiuni metodologice (relativ la inițializarea parametrilor): pr. 22.

EM pentru modelarea de mixturi de distribuții probabiliste

- varianta generală: *An Introduction to Expectation-Maximization*, Dahua Lin;
- diverse instanțe ale acestei „variante”: mixturi [de distribuții] Bernoulli (pr. 14, 4), mixturi [de distribuții] categoriale (pr. 5, 15) mixturi [de distribuții] Poisson (pr. 18), mixturi [de distribuții] Gamma (pr. 19).

Alte instanțe/aplicații ale schemei algoritmice EM

- EM pentru estimarea unui parametru [de tip probabilitate] pentru o distribuție discretă [în ocurență, o distribuție categorială], în condițiile existenței unei variabile „neobservabile”: pr. 3. Similar pentru distribuția multinomială: pr. 13;
- EM pentru estimarea tuturor parametrilor unei distribuții categoriale: pr. 12;
- EM pentru estimarea parametrului unei distribuții Poisson în condițiile în care o parte din valorile date lipsesc (pr. 10);
- EM pentru estimarea parametrilor a două distribuții probabiliste atunci când se dau instanțe care sunt generate de suma celor două distribuții: distribuții exponențiale (pr. 8), sau distribuții gaussiene (pr. 20);
- algoritmul Bayes Naiv ne-supervizat, i.e. algoritmul EM pentru [modelare] de mixturi de distribuții categoriale multi-variate, cu presupunerea de independență condițională a atributelor de intrare în raport cu atributul de ieșire (eticheta): pr. 7 (varianta de asignare “soft” a instanțelor la cluster) și pr. 17 (varianta “hard”);

- EM pentru estimarea probabilității de selecție a unei componente din cadrul unei mixturi [i.e. combinație liniară] de două distribuții probabiliste oarecare: pr. 9;
- EM pentru modelul [mixturii] domeniilor semantice (engl., topic model) pentru clusterizare de documente: pr. [6 și] 16.
- EM pentru estimarea [nu în sens MLE, cum a fost cazul până aici, ci] în sens MAP: pr. 20.

Rețele neuronale artificiale: Privire de ansamblu

A1. Noțiuni preliminare

- funcție matematică; compunere de funcții reale;
calculul valorii unei funcții pentru anumite valori specificate pentru argumentele/variabilele ei;
- funcție prag (sau, treaptă), funcție liniară, funcție sigmoidală (sau, logistică), funcție sigmoidală generalizată;
separabilitate liniară pentru o mulțime de puncte din \mathbb{R}^d ;
- ecuații asociate dreptelor în plan / planelor în spațiu / hiper-planelor în spațiul \mathbb{R}^d ;
ecuația dreptei în plan care trece prin două puncte date;
semnele asociate punctelor din semi-planele determinate de o dreaptă dată în plan;
- derivate ale funcțiilor elementare de variabilă reală; derivate parțiale
- vectori; operații cu vectori, în particular produsul scalar al vectorilor (\cdot);
- metoda gradientului descendent (ca metoda de optimizare); avantaje și dezavantaje; pr. 10, 23, 35, 36.

A2. [Câteva] noțiuni specifice

- *unități* neuronale artificiale (sau, *neuroni* artificiali, *perceptroni*);
tipuri de neuroni artificiali: neuroni-prag, liniari, sigmoidali;
componente ale unui neuron artificial: input, componenta de sumare, componenta / funcția de activare, output;
funcția matematică reprezentată / calculată de un neuron artificial;
- *rețea* neuronală artificială; rețele de tip feed-forward;
niveluri / straturi de neuroni, niveluri ascunse, niveluri de ieșire;
ponderi asociate conxiunilor dintr-o rețea neuronală artificială;
funcția matematică reprezentată / calculată de o rețea neuronală artificială;
granițe și zone de decizie determinate de o rețea neuronală artificială;
funcția de eroare / cost (engl., loss function).

A3. Câteva *proprietăți* de *expresivitate* ale rețelelor neuronale artificiale

- (P0) Toate cele trei tipuri de neuroni artificiali (prag, liniar, sigmoidal) produc *separatori liniari*.
Consecință: Conceptul XOR nu poate fi reprezentat / învățat cu astfel de „dispozitive” simple de clasificare.
- (P0') Rețelele neuronale artificiale pot determina granițe de decizie neliniare (și, în consecință, pot reprezenta concepte precum XOR).
Observație: Rețele de unități sigmoidale pot determina granițe de decizie curbilinii (vezi pr. 8).

- (P1) Rețele de neuroni diferite (ca structură și / sau tipuri de unități) pot să calculeze o aceeași funcție. (Vezi pr. 3 și pr. 1.c vs. pr. 2).
(P1') Dată o topologie de rețea neuronală (i.e., graf de unități neuronale al căror tip este lăsat nespecificat), este posibil ca plasând în noduri unități de un anumit tip să putem reprezenta / calcula o anumită funcție, iar schimbând tipul unora dintre unități (sau al tuturor unităților), funcția respectivă să nu mai potă fi calculată. (Vezi pr. 4 vs. pr. 34.³³¹)
- (P2) Orice unitate liniară situată pe un nivel ascuns poate fi „absorbită” pe nivelul următor (pr. 33).
- (P3) Orice funcție booleană poate fi reprezentată cu ajutorul unei rețele neuronale artificiale având doar două niveluri de perceptroni-prag (pr. 5).
- (P4) Orice funcție definită pe un interval mărginit din \mathbb{R} , care este continuă în sens Lipschitz, poate fi aproximată oricât de bine cu ajutorul unei rețele neuronale care are un singur nivel ascuns (pr. 7).

B. Algoritmi de antrenare a neuronilor artificiali folosind metoda gradientului descendent

- algoritmul de antrenare a unității liniare: pr. 37
vezi T. Mitchell, *Machine Learning*, p. 93, justificare: p. 91-92; convergența: p. 95; exemplu de aplicare: pr. 11
variante incrementală a algoritmului de antrenare a unității liniare: TM, ML book, p. 93-94; despre convergența acestei variante (ca aproximare a variantei precedente (“batch”)): TM, ML book, p. 93 jos;
- algoritmul de antrenare a perceptronului-prag și convergența: TM, ML book, p. 88-89; exemplu de aplicare: pr. 12;
- algoritmul de antrenare a perceptronului sigmoidal și justificarea sa teoretică: TM, ML book, p. 95-97;
- algoritmul *Perceptron* (!) al lui Rosenblatt; exemplu de aplicare: 17, 39;
- deducerea regulii de actualizare a ponderilor pentru tipuri particulare de perceptroni: pr. 13, 25.a, 38, 14.a;
- o justificare probabilistă (gen ipoteză de tip *maximum likelihood*) pentru minimizarea sumei pătratelor erorilor [la deducerea regulii de antrenare] pentru perceptronul liniar: pr 14.b;
- exemple de [folosire a unei] alte funcții de cost/pierdere/penalizare (engl., loss function) decât semi-suma pătratelor erorilor: suma costurilor de tip log-sigmoidal, pr. 15 (pentru perceptronul liniar), o funcție de tip cross-entropie, pr. 16 (pentru perceptronul sigmoidal).

B'. Perceptronul Rosenblatt și rezultate de convergență

- exemplu de aplicare [adică, învățare cu perceptronul Rosenblatt]: pr. 17.

³³¹Problemele 1.d și pr. 32 au în vedere o chestiune similară, însă pentru rețele cu topologii diferite: o anumită extensie a funcției XOR nu poate fi reprezentată pe rețele de neuroni-prag care au un singur nivel ascuns.

- câteva *proprietăți* simple ale perceptronului Rosenblatt: pr. 18.
- rezultate de convergență de tip “mistake bound” pentru [algoritmul de antrenare pentru] perceptronul-prag [în varianta] Rosenblatt: pr. 19, 40;
pentru perceptronul-prag (clasic): pr. 42;
învățare online cu perceptronul-prag de tip Rosenblatt: pr. 41;
- *Perceptronul* kernelizat [dual]: pr. 24; particularizare pentru cazul nucleului RBF: pr. 50.

C. Antrenarea rețelelor neuronale artificiale:

algoritmul de retro-propagare pentru rețelele feed-forward

- T. Mitchell, *Machine Learning*, p. 98: pseudo-cod pentru rețele cu unități de tip sigmoidal, cu 2 niveluri, dintre care unul ascuns; pentru deducerea regulilor de actualizare a ponderilor (în cazul mai general rețelelor al feed-forward (de unități sigmoidale) cu oricâte niveluri, vezi p. 101-103);
pr. 20: deducerea regulilor de actualizare a ponderilor în cazul rețelelor cu 2 niveluri, având însă unități cu funcție de activare oarecare (derivabilă);
- aplicare: pr. 21, 43, 44;
- prevenirea overfitting-ului:
folosirea unei componente de tip „moment” în expresia regulilor de actualizare a ponderilor: pr. 46;
regularizare: introducerea unei componente suplimentare în funcția de optimizat: pr. 22;
- cazul folosirii unei funcții de activare de tip tangentă hiperbolică: pr. 45;
- cazul folosirii unei funcții de cost/penalizare/eroare de tip cross-entropie: pr. 48;
- fenomenul de „dispariție” a gradientului [în cazul aplicării algoritmului de retro-propagare] pentru rețele neuronale profunde (engl., deep neural networks) care folosesc funcția de activare sigmoidală: pr. 27;
execuția manuală a unei iterații a algoritmului de retro-propagare în cazul unei rețele neuronale simple, având un singur nivel ascuns, cu unități ce folosesc funcția de activare ReL: pr. 49.

D. Rețele neuronale profunde:

chestiuni introductive: pr. 28 și pr. 55.