

# Exerciții de învățare automată

Liviu Ciortuz, Alina Munteanu, Elena Bădăraș

2019

ISBN: 978-973-0-30467-1

## 1. Fundamente:

probabilități, variabile aleatoare, teoria informației,  
funcții-nucleu, și metode de optimizare

### Sumar

#### Evenimente aleatoare și formula lui Bayes

- *funcția de probabilitate* – câteva proprietăți care derivă din definiția ei: ex. 1, ex. 64;
- calcularea unor *probabilități* elementare: ex. 2.a, ex. 61;
- calcularea unor *probabilități condiționate*: ex. 2.b, ex. 3, ex. 62.a;
- *regula de multiplicare*: ex. 64.b;
- *formula probabilității totale*: ex. 64.e;  
formula probabilității totale – varianta condițională: ex. 67.cd;
- *independența* evenimentelor aleatoare: ex. 4, ex. 5, ex. 62.bc;
- *independența condițională* a evenimentelor aleatoare – legătura dintre forma „tare” și forma „slabă” a definiției pentru această noțiune: ex. 63;
- *formula lui Bayes* – aplicare: ex. 6, ex. 7, ex. 65, ex. 66;  
formula lui Bayes – varianta condițională: ex. 67.b;
- recapitulare: probabilități elementare și probabilități condiționate – câteva proprietăți ( $A/F$ ): ex. 8, ex. 67.

#### Variabile aleatoare

- funcție de *distribuție cumulativă* (engl., cumulative distribution function, c.d.f.), un exemplu: ex. 27;
- proprietatea de *liniaritate* a *mediei*: ex. 9.a;
- *varianța* și *covarianța*: proprietăți de tip *caracterizare*: ex. 9.bc;
- covarianța oricăror două variabile aleatoare independente este 0: ex. 10;  
reciproca acestei afirmații nu este adevărată în general: ex. 11.a; totuși ea are loc dacă variabilele sunt de tip binar (adică iau doar valorile 0 și 1): ex. 11.b, ex. 71;
- o *condiție suficientă* pentru ca  $\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$ : independența variabilelor  $X$  și  $Y$ : ex. 19.c
- pentru *variabile aleatoare discrete*:
- calcularea mediilor și a varianțelor – exemplificare: ex. 12, ex. 68, ex. 70.a;

- definirea unei *variabile-indicator* cu ajutorul unui eveniment aleator; calcularea mediei acestei variabile: ex. 69;
- regula de *înlanțuire* (pt. var. aleat.) – aplicare: ex. 13;
- regula de *multiplicare* (pt. var. aleat.), varianta condițională – demonstrație: ex. 14;
- *independență*: ex. 70.b, ex. 72.a, ex. 73.a;  
*independență condițională*: ex. 15, ex. 72.b, ex. 73.b, ex. 74.b;  
o condiție suficientă pentru independența condițională: ex. 75;
- pentru *variabile aleatoare continue*:
  - dată fiind o funcție care depinde de un parametru real, să se calculeze valoarea respectivului parametru astfel încât funcția respectivă să fie o *funcție de densitate de probabilitate* (engl., probability density function, p.d.f.): ex. 16.a, ex. 17, ex. 76;
  - dat fiind un p.d.f., să se calculeze o anumită probabilitate: ex. 16.c, ex. 77;
  - variabile aleatoare discrete vs. variabile aleatoare continue; p.m.f. vs. p.d.f.: ex. 78;
  - variabile aleatoare discrete și variabile aleatoare continue; independența și calculul de medii: ex. 79;
- *vector de variabile aleatoare*:
  - o proprietate: matricea de covarianță este simetrică și pozitiv definită: ex. 18;
  - calculul matricei de covarianță când asupra vectorului de variabile aleatoare operăm transformări liniare: ex. 80;
- recapitulare: ex. 19, ex. 81.

## Distribuții probabiliste uzuale

- distribuții discrete
  - distribuția *Bernoulli*:  
suma de variabile identic distribuite; media sumei: ex. 20;  
*mixtură* de distribuții Bernoulli: ex. 85;
  - distribuția *binomială*:  
verificarea condițiilor de definiție pentru p.m.f.: ex. 21.a;  
calculul mediei și al varianței: ex. 21.bc;  
calcularea unor probabilități (simple și respectiv condiționale): ex. 82;
  - distribuția *categorială*:  
calcularea unor probabilități și a unor medii: ex. 83;  
*mixtură* de distribuții categoriale: calculul mediei și al varianței: ex. 83;
  - distribuția *geometrică*:  
calcularea numărului „așteptat” / mediu de „observații” necesare pentru ca un anumit eveniment să se producă: ex. 84;
  - distribuția *Poisson*:  
verificarea condițiilor de definiție pentru p.m.f., calculul mediei și al varianței: ex. 22;
- distribuții continue

- distribuția *continuuă uniformă*:  
exemplu de distribuție continuă uniformă uni-variată; calcularea mediei și a varianței: ex. 87;  
calculul unei p.d.f. corelate, pornind de la două variabile [urmând distribuții continue uniforme uni-variate] independente; calculul unei anumite probabilități, folosind p.d.f. corelată: ex. 24;  
exemplu de distribuție continuă uniformă bivariată; calcularea unei p.d.f. condiționale; verificarea independenței celor două distribuții marginale; calculul mediilor unor p.d.f.-uri condiționale: ex. 86, ex. 79;
- distribuția *exponențială*:  
verificarea condițiilor de definiție pentru p.d.f.,  
calculul mediei și al varianței: ex. 25.a;
- distribuția *Gamma*:  
verificarea condițiilor de definiție pentru p.d.f., calculul mediei și al varianței: ex. 25.b;
- distribuția *gaussiană uni-variată*:  
verificarea condițiilor de definiție pentru p.d.f., calculul mediei și al varianței: ex. 26;  
„standardizare” (i.e., reducerea cazului nestandard la cazul standard): ex. 27;
- distribuția *gaussiană multi-variată*:  
matrice simetrice și pozitiv definite<sup>1</sup> o proprietate de tip *factorizare* folosind *vectori proprii*: ex. 29;  
p.d.f.-ul distribuției gaussiene multi-variate este într-adevăr p.d.f.:<sup>2</sup> ex. 30;  
o *proprietate* importantă, în cazul în care matricea de covarianță este diagonală: p.d.f.-ul corelat este produsul p.d.f.-urilor marginale (care sunt independente): ex. 28;  
distribuția *gaussiană bivariată*: exemplificare; calcularea explicită a p.d.f.-ului, dat fiind vectorul de medii și matricea de covarianță: ex. 88;  
o *proprietate* pentru distribuția *gaussiană bivariată*: distribuția condițională a unei componente în raport cu cealaltă componentă este tot de tip gaussian; calculul parametrilor acestei distribuții condiționale: ex. 31;
- *mixturi* de distribuții *gaussiene multi-variate*:  
exprimarea vectorului de medii și a matricei de covarianță în funcție de mediile și matricele de covarianță ale distribuțiilor componente: ex. 89;
- *mixturi* de distribuții oarecare:  
calculul mediilor și al varianțelor (în funcție de mediile și matricele de covarianță ale distribuțiilor componente): ex. 90;
- distribuția Bernoulli și distribuția normală standard:  
*intervale de încredere*, *teorema limită centrală*; aplicație la calculul erorii reale a unui clasificator: ex. 32;
- chestiuni recapitulative (corespondența dintre nume de distribuții și expresiile unor p.d.f.-uri date): ex. 91.

---

<sup>1</sup> Așa sunt matricele de covarianță ale variabilelor gaussiene multi-variate.

<sup>2</sup> Adică satisface condițiile din definiția noțiunii de p.d.f.

## Elemente de teoria informației

- re-descoperirea definiției entropiei, pornind de la un set de proprietăți dezirabile: ex. 33;
- *definiții și proprietăți imediate* pentru entropie, entropie corelată, entropie condițională specifică, entropie condițională medie, câștig de informație: ex. 34;  
*exemplificarea* acestor noțiuni (varianta discretă): ex. 35, ex. 36, ex. 92;  
un *exemplu* de calcul pentru entropia unei variabile continue (distribuția exponențială): ex. 37;  
o *proprietate* a entropiei: ne-negativitatea: ex. 34.a, ex. 98.a;  
o *margină superioară* pentru valoarea entropiei unei variabile aleatoare discrete: ex. 93;  
o *proprietate* a câștigului de informație: ne-negativitatea: ex. 38.c, ex. 96;
- o *aplicație* pentru câștigul de informație: *selecția de trăsături*: ex. 39;
- *entropia corelată*: forma particulară a relației de „înlănțuire” în cazul variabilelor aleatoare independente: ex. 40, ex. 94, ex. 98.c;  
Entropie corelată și condiționată: formula de „înlănțuire” condițională: ex. 95;
- *entropia relativă*: definiție și proprietăți elementare; exprimarea câștigului de informație cu ajutorul entropiei relative: ex. 38;
- *cross-entropie*: definiție, o proprietate (ne-negativitatea) și un exemplu simplu de calculare a valorii cross-entropiei: ex. 41;  
un exemplu de *aplicație* pentru cross-entropie: selecția modelelor probabiliste: ex. 97;
- *inegalitatea lui Gibbs*: un caz particular; comparație între valorile entropiei și ale cross-entropiei: ex. 42.

## Funcții-nucleu

- aflarea funcției de „mapare” a trăsăturilor care corespunde unei funcții-nucleu date: ex. 43, ex. 99, ex. 100.a;  
comparații asupra numărului de operații efectuate la calcularea valorii unor funcții-nucleu (în spațiul inițial vs. spațiul nou de „trăsături”): ex. 100.b;  
calculul distanțelor euclidiene în spații de „trăsături” folosind doar funcții-nucleu: ex. 47;
- *teorema lui Mercer* (1909): condiții necesare și suficiente pentru ca o funcție să fie funcție-nucleu: ex. 44.ab;
- rezultate de tip „constructiv” pentru [obținerea de noi] funcții-nucleu: ex. 44.c, 45, ex. 46, ex. 101, 103, ex. 52.b;  
„normalizarea” funcțiilor-nucleu: ex. 102;
- o inegalitate [derivată din inegalitatea Cauchy-Buniakovski-Schwarz], care furnizează o margine superioară pentru  $K(x, x')$ , valoarea absolută a unei funcții-nucleu oarecare: ex. 104;
- un exemplu de funcție-nucleu care servește la a măsura similaritatea dintre două imagini oarecare: ex. 48;
- o funcție-nucleu particulară, care asigură separabilitate liniară [în spațiul de trăsături] pentru orice set de instanțe de antrenament: ex. 105;

- funcția-nucleu gaussiană / *funcția cu baza radială* (engl., Radial Basis Function, RBF):
  - demonstrația faptului că RBF este într-adevăr funcție-nucleu: ex. 49;
  - funcția de „mapare” corespunzătoare funcției-nucleu RBF la valori într-un spațiu [de „trăsături”] de dimensiune infinită: ex. 50;
  - proprietăți simple ale nucleului RBF: ex. 51, ex. 106;
  - orice mulțime de instanțe distincte, având orice etichetare posibilă, este separabilă liniar în spațiul de „trăsături” dacă se folosește nucleul RBF cu parametrul ales în mod convenabil: ex. 26.a de la capitolul *Mașini cu vectori-suport*;
- recapitulare (A/F): ex. 52, ex. 108.

## Metode de optimizare în învățarea automată

- (P0) definiții, caracterizări și câteva proprietăți pentru funcții convexe: ex. 53;
- *metoda gradientului*, *metoda lui Newton*; exemplificare: ex. 54;
  - (P1) condiții suficiente pentru convergența metodei gradientului: ex. 110;
  - (P2) o proprietate interesantă a metodei lui Newton: în cazul oricărei funcții de gradul al doilea (de una sau mai multe variabile), aplicarea acestei metode de optimizare implică / necesită execuția unei singure iterații: ex. 111;
  - (P3) *reparametrizarea liniară* a atributelor nu afectează [rezultatele obținute cu] metoda lui Newton, însă afectează metoda gradientului: ex. 112;
- metoda dualității Lagrange:
  - (P4) demonstrarea proprietății de *dualitate slabă*: ex. 55;
  - (P5) demonstrarea unei părți din *teorema KKT*: ex. 56;
 exemple de aplicare: ex. 57, ex. 58, ex. 113;  
 un exemplu de problemă de optimizare convexă pentru care condițiile KKT nu sunt satisfăcute: ex. 114;
- două variante a algoritmului Perceptron,<sup>3</sup> pentru care relația de actualizare a ponderilor se obține rezolvând [câte] o problemă de optimizare [convexă] cu restricții;
- chestiuni recapitulative: ex. 109.

---

<sup>3</sup>Vedeți problema 16 de la capitolul *Rețele neuronale artificiale*.

## 2. Metode de estimare a parametrilor; metode de regresie

### Sumar

#### Noțiuni preliminare

- elemente de calcul vectorial (în particular, produsul scalar) și de calcul matriceal: ex. 29 de la cap. *Fundamente*; norma  $L_2$  (euclidiană) și norma  $L_1$ : ex. 14, ex. 47, ex. 49; calculul derivatelor parțiale [de ordinul întâi și al doilea]: ex. 17; reguli de derivare cu argumente vectoriale: ex. 12;
- metode de optimizare (în speță pentru aflarea maximumului / minimumului unei funcții reale, derivabile): metoda analitică, metoda gradientului, metoda lui Newton; exemplificare: ex. 54 de la capitolul de Fundamente;

#### Estimarea parametrilor unor distribuții probabiliste uzuale

- distribuția Bernoulli: ex. 1 (+MAP, folosind distr. Beta), ex. 2, ex. 30 (un caz particular), ex. 29 (bias-ul și varianța estimatorului MLE);
- distribuția categorială: ex. 31, ex. 32 (+MAP, folosind distr. Dirichlet);
- distribuția geometrică: ex. 33 (+MAP, folosind distr. Beta);
- distribuția Poisson: ex. 3 (+MAP, folosind distr. Gamma);
- distribuția uniformă continuă: calcul de probabilități și MLE: în  $\mathbb{R}$ : ex. 4, ex. 34, ex. 35; în  $\mathbb{R}^2$ : ex. 5, ex. 36;
- distribuția gaussiană uni-variată:
  - MLE pt.  $\mu$ , considerând  $\sigma^2$  cunoscut: ex. 2 (+MAP, folosind distribuția gaussiană), ex. 38 (+analiza discriminativă);
  - MLE pt.  $\sigma^2$ , atunci când nu se impun restricții asupra lui  $\mu$ : ex. 8 (+deplasare);
  - MLE pt.  $\sigma^2$ , atunci când  $\mu = 0$ : ex. 39 (+nedeplasare);
- distribuția gaussiană multi-variată: ex. 10, ex. 40 (+MAP, folosind distr. Gauss-Wishart);
- distribuția exponențială: ex. 6, ex. 37 (+MAP, folosind distr. Gamma);
- distribuția Gamma: ex. 9 și ex. 41 (ultimul, folosind metoda gradientului și metoda lui Newton).
- existența și unicitatea MLE: ex. 11.

#### Regresia liniară

- prezentarea *generală* a metodei regresiei liniare:<sup>4</sup>
  - MLE și corespondența cu estimarea în sens LSE (least squared errors): ex. 14.A;
  - particularizare pentru cazul uni-variat: ex. 12.ab, ex. 42;
  - exemplificare pentru cazul uni-variat (ex. 13, ex. 43) și pentru cazul bivariat (ex. 45.a, ex. 53);

---

<sup>4</sup>În mod implicit, în această secțiune se va considera că termenul-zgomot este modelat cu distribuția gaussiană (dacă nu se specifică altfel, în mod explicit).

- (P1) *scalarea atributelor* nu schimbă predicțiile obținute (pentru instanțele de test) cu ajutorul *formulelor analitice*: ex. 15, 59.a;
- (P2) adăugarea de noi trăsături / attribute nu mărește suma pătratelor erorilor: ex. 48;
- o *proprietate surprinzătoare* a regresiei liniare: adăugarea câtorva „observații” suplimentare poate conduce la modificarea radicală a valorilor optime ale parametrilor de regresie: CMU, 2014 fall, Z. Bar-Joseph, W. Cohen, HW2, pr. 4;
- [rezolvarea problemei de] regresie liniară folosind *metoda lui Newton*: ex. 17;
- MAP și corespondența cu *regularizarea* de normă  $L_2$  (regresia *ridge*): ex. 14.C; particularizare pentru cazul uni-variabil: ex. 12.c;
- *regularizarea* de normă  $L_1$  (regresia *Lasso*): ex. 47.a;
- (P3) efectul de diminuare a ponderilor (engl., *weight decay*) în cazul regularizării de normă  $L_2$  (respectiv  $L_1$ ) a regresiei liniare, în comparație cu cazul neregularizat: ex. 47.b;
- bias-ul și [co]varianța estimatorului regresiei liniare; bias-ul regresiei *ridge*: ex. 16;
- regresia polinomială [LC: mai general: folosirea așa-numitelor *funcții de bază*]: ex. 14.B; exemplificare pentru cazul bivariat: CMU, 2015 spring, T. Mitchell, N. Balcan, HW4, pr. 1;
- cazul regresiei liniare cu termen de regularizare  $L_2$  (regresia *ridge*): deducerea regulilor de actualizare pentru *metoda gradientului ascendent*: varianta “batch” / “steepest descent”: ex. 18.a; și varianta stohastică / secvențială / “online”: ex. 18.b; exemplu de aplicare: ex. 46;
- cazul regresiei liniare cu termen de regularizare  $L_1$  (regresia *Lasso*): rezolvare cu *metoda descenderii pe coordonate* (engl., “coordinate descent”): ex. 49; rezolvare cu *metoda sub-gradientului* (aplicare la selecția de trăsături): CMU, 2009 fall, C. Guestrin, HW2, pr. 2;
- regresia liniară în cazul zgomotului modelat cu distribuția *Laplace* (în locul zgomotului gaussian): ex. 19.B; exemplificare pentru cazul bivariat: ex. 45.c; rezolvare în cazul uni-variabil [chiar particularizat] cu ajutorul derivatei, acolo unde aceasta există: ex. 50;
- regresia liniară și *overfitting-ul*: ex. 22;
- regresie liniară folosită pentru *clasificare*: exemplificare: ex. 53;
- cazul *multi-valuat* al regresiei liniare, reducerea la cazul uninomial: ex. 52;
- regresia liniară cu regularizare  $L_2$  (regresia *ridge*), *kernel-izarea* ecuațiilor „normale”: ex. 20; (P4) folosind nucleu RBF, eroarea la antrenare devine 0 atunci când parametrul de regularizare  $\lambda$  tinde la 0: ex. 21;
- *regresia liniară ponderată*: ex. 19.A; particularizare / exemplificare pentru cazul bivariat: ex. 45.b;

- o *proprietate* a regresiei liniare local-ponderate [demonstrată în cazul uni-variat]: „netezirea“ liniară: ex. 51; cazul multi-valuat, cu regularizare  $L_2$ : Stanford, 2015 fall, Andrew Ng, midterm, pr. 2;
- o *regresia liniară (kernelizată) local-ponderată, ne-parametrică*: particularizare / exemplificare pentru cazul uni-variat, cu nucleu gaussian: CMU, 2010 fall, Aarti Singh, midterm, pr. 4.

## Regresia logistică

- prezentare generală,
  - (•) calculul funcției de log-verosimilitate, estimarea parametrilor în sens MLE, folosind *metoda gradientului* (i.e., deducerea regulilor de actualizare a parametrilor): ex. 23, 59.b;
  - particularizare* pentru cazul datelor din  $\mathbb{R}^2$ : ex. 54 (inclusiv *regularizare*  $L_1$  / estimarea parametrilor în sens MAP, folosind o distribuție a priori Laplace);
- (P0) *granița de decizie* pentru regresia logistică: ex. 54.d;
- (P1) funcția de log-verosimilitate în cazul regresiei logistice este concavă (decă are un maxim global), fiindcă matricea hessiană este pozitiv definită: ex. 24;
 

*Observație:* Demonstrația furnizează tot ce este necesar pentru obținerea [ulterioară a] relației de actualizare a parametrilor la aplicarea *metodei lui Newton* în cazul regresiei logistice;
- (P2) analiza efectului duplicării atributelor: ex. 55;
- (P3) efectul de diminuare a ponderilor (engl., weight decay) în cazul regularizării de normă  $L_2$  a regresiei logistice — adică la estimarea parametrilor în sens MAP, folosind ca distribuție a priori distribuția gaussiană multi-variată sferică —, în comparație cu cazul estimării parametrilor în sensul MLE: ex. 25;
- *Variante / extensii* ale regresiei logistice:
  - regresia logistică local-ponderată, cu regularizare  $L_2$ :
    - (•) calcularea vectorului gradient și a matricei hessiene (necesare pentru aplicarea metodei lui Newton în acest caz): ex. 56;
  - regresia logistică kernel-izată:
    - (•) adaptarea metodei gradientului: ex. 26;
  - regresia logistică  $n$ -ară (așa-numita regresie *softmax*), cu regularizare  $L_2$ :
    - (•) calculul funcției de log-verosimilitate, deducerea regulilor de actualizare a ponderilor, folosind metoda gradientului: ex. 27;
- (P4) o [interesantă] proprietate comună pentru regresia liniară și regresia logistică: ex. 57;
- întrebări (cu răspuns A/F) cu privire la aplicarea *metodei lui Newton* comparativ cu *metoda gradientului* (în contextul rezolvării problemelor de regresie liniară și / sau regresie logistică): ex. 59.c;
- comparații între regresia logistică și alți clasificatori (Bayes Naiv, ID3): ex. 54.c, ex. 58.ab.



### 3. Arbori de decizie

#### Sumar

#### Noțiuni preliminare

- partiție a unei mulțimi: ex. 64 de la cap. *Fundamente*;
- proprietăți elementare ale funcției logaritm; formule uzuale pentru calcule cu logaritmi;
- entropie, definiție: T. Mitchell, *Machine Learning*, 1997 (desemnată în continuare simplu prin *cartea ML*), pag. 57; ex. 2.a, ex. 37.a;
- entropie condițională specifică: ex. 15.a;
- entropie condițională medie: ex. 2.cd;
- câștig de informație (definiție: *cartea ML*, pag. 58): ex. 2.cd, ex. 5.a, ex. 31, ex. 37.b;
- *arbori de decizie*, văzuți ca *structură de date*: ex. 1, ex. 29  
și, respectiv, ca program în logica propozițiilor: ex. 2.e, ex. 36.bc;
  - (P0) *expresivitatea arborilor de decizie* cu privire la *funcții boolene*: ex. 30;
- *spațiu de versiuni* pentru un concept (de învățat): ex. 1, ex. 29, ex. 35.

#### Algoritmul ID3

- pseudo-cod: *cartea ML*, pag. 56;
- *bias-ul inductiv*: *ibidem*, pag. 63-64;
- exemple simple de aplicare: ex. 2, ex. 3, ex. 5, ex. 34, ex. 35, ex. 37, ex. 38;
- ID3 ca algoritm *per se*:
  - este un *algoritm de căutare*;  
*spațiul de căutare* — mulțimea tuturor arborilor de decizie care se pot construi cu atributele de intrare în nodurile de test și cu valorile atributului de ieșire în nodurile de decizie — este de dimensiune exponențială în raport cu numărul de atribute: ex. 1, ex. 3, ex. 29, ex. 35;  
ID3 are ca *obiectiv* căutarea unui arbore / *model* care *i.* să explice cât mai bine datele (în particular, atunci când datele sunt *consistente*, modelul trebuie să fie *consistent* cu acestea), *ii.* să fie cât mai *compact*, din motive de *eficiență* la generalizare și *iii.* în final să aibă o [cât mai] bună putere de *generalizare*;<sup>5</sup>
  - ID3 ar putea fi văzut și ca algoritm de *optimizare*;<sup>6</sup>
  - *greedy*  $\Rightarrow$  nu garantează obținerea soluției optime d.p.v. al numărului de niveluri / noduri:  
ex. 4, ex. 22.a, ex. 36 (vs. ex. 35.b, ex. 3.b), ex. 44;

<sup>5</sup>LC: Alternativ, putem spune că algoritmul ID3 produce o *structură* de tip *ierarhie* (arbore) între diferite *partiționări* ale setului de instanțe de antrenament, această ierarhie fiind generată pe baza *corespondenței* dintre atributul de *ieșire* și atributele de *intrare*, care sunt adăugate la model câte unul pe rând.

<sup>6</sup>LC: Am putea să-l interpretăm pe ID3 ca fiind un algoritm care caută între diferitele *distribuții probabiliste* discrete care pot fi definite pe setul de date de antrenament una care să satisfacă cerința de *ierarhizare* (vedeți mai jos), și pentru care entropia să fie *minimală* (vedeți proprietatea de structuralitate de la ex. 33 de la cap. *Fundamente*. Cerința ca arborele ID3 să fie *minimal* (ca număr de niveluri / noduri) este însă mai importantă, mai practică și mai ușor de înțeles.

- de tip *divide-et-impera* ( $\Rightarrow$  “*Iterative Dichotomizer*”), recursiv;
    - *1-step look-ahead*;
  - complexitate de timp (vedeți 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, iar  $m$  este numărul de exemple;
- ID3 ca algoritm de învățare automată:
- *bias*-ul inductiv al algoritmului ID3:  
[dorim ca modelul să aibă structură ierarhică, să fie compatibil / consistent cu datele dacă acestea sunt consistente (adică, necontradictorii), iar *arboarele* [produs de] *ID3* trebuie să aibă un număr cât mai mic de niveluri / noduri;
  - algoritm de învățare de tip “*eager*”;
  - *analiza erorilor*:  
la antrenare: ex. 7.a, ex. 10.a, ex. 40;  
[acuratețe la antrenare: ex. 6;]  
la validare: CMU, 2003 fall, T. Mitchell, A. Moore, midterm, pr. 1;  
la  $n$ -fold cross-validare  
la cross-validare leave-one-out (CVLOO): ex. 10.b, ex. 43.bc;
  - *robustețe la „zgomot” și overfitting*: ex. 10, ex. 22.bc, ex. 43, ex. 66.b;
  - *zone de decizie și granițe de separare / decizie* pentru arbori de decizie cu variabile continue: ex. 10, ex. 42, ex. 43, ex. 44.

## Extensii / variante ale algoritmului ID3

- *attribute cu valori continue*: ex. 10-12, ex. 15.c, ex. 41-45; cap. *Învățare bazată pe memorare*, ex. 11.b;  
alte variante de partiționare a intervalelor de valori pentru atributele continue: ex. 47;
- *attribute discrete cu multe valori*: ex. 14;
- *attribute cu valori nespecificate / absente pentru unele instanțe*;
- *attribute cu diferite costuri asociate*: ex. 15;
- *reducerea caracterului “eager” al învățării*: ex. 17;
- *reducerea caracterului “greedy” al învățării*:  
IG cu “2-step look-ahead”: ex. 18, ex. 19;  
variante de tip “look-ahead” specifice atributelor continue: ex. 48;  
3-way splitting (sau, mai general,  $n$ -way splitting) pentru atributele continue: ex. 13, ex. 46;
- *folosirea altor măsuri de „impuritate” în locul câștigului de informație*:  
Gini Impurity, Misclassification Impurity: ex. 16;
- *reducerea overfitting-ului*:  
reduced-error pruning (folosind un set de date de validare):  
cartea ML, pag. 69-71; A. Cornuéjols, L. Miclet, 2nd ed., pag. 418-421;  
rule post-pruning: cartea ML, pag. 71-72; ex. 50  
top-down vs. bottom-up pruning: ex. 20, ex. 49;  
pruning folosind testul statistic  $\chi^2$ : ex. 21, ex. 51.

## Proprietăți ale arborilor ID3

- (P1) arborele produs de algoritmul ID3 este *consistent* (adică, în concordanță) cu datele de antrenament, dacă acestea sunt *consistente* (adică, necontradictorii). Altfel spus, *eroarea la antrenare* produsă de algoritmul ID3 pe orice set de *date consistente* este 0: ex. 2-4, ex. 35-36;
- (P2) arborele produs de algoritmul ID3 *nu* este în mod neapărat *unic*: ex. 3, ex. 35;
- (P3) arborele ID3 *nu* este neapărat *optimal* (ca nr. de noduri / niveluri): ex. 4, ex. 22, ex. 36;
- (P4) influența *atributelor identice* și, respectiv, a *instanțelor multiple* asupra arborelui ID3: ex. 8;
- (P5) o *margină superioară* pentru *eroarea la antrenare* a algoritmului ID3, în funcție de numărul de valori ale variabilei de ieșire): ex. 7.b;
- (P6) o aproximare simplă a numărului de *instanțe greșit clasificate* din totalul de  $M$  instanțe care au fost asignate la un nod frunză al unui arbore ID3, cu ajutorul entropiei ( $H$ ) nodului respectiv: ex. 39;
- (P7) *granițele de separare / decizie* pentru arborii ID3 cu atribute de intrare continue sunt întotdeauna paralele cu axele de coordonate: ex. 10, ex. 12, ex. 42, ex. 43, ex. 44 și cap. *Învățare bazată pe memorare*, ex. 11.b.  
*Observație:* Următoarele trei proprietăți se referă la arbori de decizie în general, nu doar la arbori ID3.
- (P8) *adâncimea maximă* a unui arbore de decizie, când atributele de intrare sunt categoricale: numărul de atribute: ex. 52.c;
- (P9) o *margină superioară* pentru *adâncimea* unui arbore de decizie când atributele de intrare sunt continue, iar datele de antrenament sunt (ne)separabile liniar: ex. 11;
- (P10) o *margină superioară* pentru numărul de *noduri-frunză* dintr-un arbore de decizie, în funcție de numărul de exemple și de numărul de atribute de intrare, atunci când acestea (atributele de intrare) sunt binare: ex. 9;

## Învățare automată de tip ansamblist folosind arbori de decizie: Algoritmul AdaBoost

- Noțiuni preliminare:  
distribuție de probabilitate discretă, factor de normalizare pentru o distribuție de probabilitate, ipoteze „slabe” (engl., weak hypothesis), compas de decizie (engl., decision stump), prag de separare (engl., threshold split) pentru un compas de decizie, prag exterior de separare (engl., outside threshold split), eroare ponderată la antrenare (engl., weighted training error), vot majoritar ponderat (engl., weighted majority vote), overfitting, ansambluri de clasificatori (vedeți ex. 64), funcții de cost / pierdere (engl., loss function) (vedeți ex. 63 și ex. 23);
- pseudo-codul algoritmului AdaBoost + convergența erorii la antrenare: ex. 23;
- exemple de aplicare: ex. 24, 54, 55, 56, 57;

- AdaBoost ca algoritm *per se*:  
algoritm iterativ, algoritm de căutare (spațiul de căutare este mulțimea combinațiilor liniare care se pot construi peste clasa de ipoteze „slabe” considerate), algoritm de optimizare secvențială (minimizează o margine superioară pentru eroarea la antrenare), algoritm greedy.
- *învățabilitate empirică  $\gamma$ -slabă*:  
definiție: ex. 23.f  
exemplificarea unor cazuri când nu există *garanție* pentru învățabilitate  $\gamma$ -slabă: ex. 59, 60;
- AdaBoost ca algoritm de *optimizare secvențială* în raport cu funcția de cost / „pierdere” negativ-exponențială: ex. 25;
- marginea de votare: ex. 26;
- *selectarea trăsăturilor* folosind AdaBoost; aplicare la clasificarea de documente: ex. 61;
- o variantă generalizată a algoritmului AdaBoost: ex. 63;
- recapitulare (întrebări cu răspuns *adevărat* / *fals*): ex. 28 și 65;
- Proprietăți ale algoritmului AdaBoost:
  - (P0) AdaBoost poate produce rezultate diferite atunci când are posibilitatea să aleagă între două sau mai multe [cele mai bune] ipoteze „slabe”: ex. 24, 54;
  - (P1)  $err_{D_{t+1}}(h_t) = \frac{1}{2}$  (ex. 23.a);  
ca o *consecință*, rezultă că ipoteza  $h_t$  nu poate fi reselectată și la iterația  $t + 1$ ; ea poate fi reselectată la o iterație ulterioară;
  - (P2) Din relația de definiție pentru distribuția  $D_{t+1}$  rezultă  $Z_t = e^{-\alpha_t} \cdot (1 - \varepsilon_t) + e^{\alpha_t} \cdot \varepsilon_t = 2\sqrt{\varepsilon_t(1 - \varepsilon_t)}$  (ex. 23.a);  
 $\varepsilon_t \in (0, 1/2) \Rightarrow Z_t \in (0, 1)$ .
  - (P3)  $D_{t+1}(i) = \frac{1}{m \prod_{t'=1}^t Z_{t'}} e^{-y_i f_t(x_i)}$ , unde  $f_t(x_i) \stackrel{\text{def.}}{=} -\sum_{t'=1}^t \alpha_{t'} h_{t'}(x_i)$  (ex. 23.b).  
Produsul  $y_i f_t(x_i)$  se numește *margine algebrică*;
  - (P4)  $err_S(H_t) \leq \prod_{t'=1}^t Z_{t'}$ , adică eroarea la antrenare comisă de ipoteza combinată produsă de AdaBoost este majorată de produsul factorilor de normalizare (ex. 23.c);
  - (P5) AdaBoost nu optimizează în mod direct  $err_S(H_t)$ , ci marginea sa superioară,  $\prod_{t'=1}^t Z_{t'}$ ; optimizarea se face în mod secvențial (greedy): la iterația  $t$  se minimizează valoarea lui  $Z_t$  ca funcție de  $\alpha_t$ , ceea ce conduce la  $\alpha_t = \ln \sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}}$  (ex. 23.d);
  - (P5')  $\varepsilon_i > \varepsilon_j \Rightarrow \alpha_i < \alpha_j$  (consecință imediată din (P5)): ex. 60.c;
  - (P6) O consecință din relația (118) și (P5):  $D_{t+1}(i) = \begin{cases} \frac{1}{2\varepsilon_t} D_t(i), & i \in M \\ \frac{1}{2(1 - \varepsilon_t)} D_t(i), & i \in C \end{cases}$
  - (P7)  $err_S(H_t)$  nu neapărat descrește de la o iterație la alta; în schimb, descreșc marginile sale superioare:  $\prod_{t'=1}^t Z_{t'}$  și  $\exp(-\sum_{t'=1}^t \gamma_{t'}^2)$  (ex. 23.e);
  - (P8) O condiție suficientă pentru învățabilitate  $\gamma$ -slabă, bazată pe marginea de votare: marginea de votare a oricărei instanțe de antrenament să fie de cel

puțin  $2\gamma$ , la orice iterație a algoritmului AdaBoost (ex. 27);

(P9) Orice mulțime formată din  $m$  instanțe din  $\mathbb{R}$  care sunt etichetate în mod consistent poate fi corect clasificată de către o combinație liniară formată din cel mult  $2m$  compași de decizie (ex. 64.a);

(P10) Orice mulțime de instanțe distincte [și etichetate] din  $\mathbb{R}$  este  $\gamma$ -slab învățabilă cu ajutorul compașilor de decizie (ex. 62).

- Alte metode de învățare ansamblistă bazate pe arbori de decizie: Bagging, Random Forests;
- Alte metode de învățare automată bazate pe arbori: arbori de regresie (CART).

## 4. Clasificare bayesiană

### Sumar

#### Noțiuni preliminare

- probabilități și probabilități condiționate;
- formula lui Bayes: ex. 5.b;  
cap. *Fundamente*, ex. 6, ex. 7, ex. 65, ex. 66;
- independența [condițională a] evenimentelor aleatoare:  
cap. *Fundamente*, ex. 4, ex. 62, ex. 63;
- independența [condițională a] variabilelor aleatoare: ex. 9, ex. 10, ex. 12, ex. 28-35; vedeți și cap. *Fundamente*, ex. 15, ex. 24, ex. 70.b, ex. 79, ex. 86;
- distribuții probabiliste corelate, marginale și condiționale: ex. 8, ex. 10, ex. 12, ex. 28; vedeți și cap. *Fundamente*, ex. 13, ex. 14;
- distribuția gaussiană: de la cap. *Fundamente*, ex. 26, ex. 27 (pentru cazul uni-variat), ex. 88 (pentru cazul bivariat), ex. 18, ex. 28, ex. 29, ex. 30 (pentru cazul multi-variat);
- estimarea parametrilor pentru distribuții de tip Bernoulli, categorial și gaussian (ultimul doar pentru cazul clasificării bayesiene de tip gaussian);<sup>7</sup>
- ipoteze MAP vs. ipoteze ML:  
formulare [ca soluții la] probleme de optimizare:<sup>8</sup> ex. 22;  
exemplificare: ex. 1, ex. 2, ex. 3, ex. 21, ex. 34;  
exemplificare în cazul arborilor de decizie: ex. 4;
- regresia logistică, chestiuni introductive:<sup>9</sup> de la cap. *Estimarea parametrilor; metode de regresie*, ex. 23.

#### Algoritmi de clasificare bayesiană

- Algoritmul Bayes Naiv și algoritmul Bayes Corelat:<sup>10</sup>  
formulare ca probleme de optimizare / estimare în sens MAP: cartea ML, pag. 167;  
pseudo-cod: vedeți slide-uri;  
exemple de aplicare: ex. 5, ex. 7, ex. 8, ex. 9, ex. 23, ex. 24, ex. 25;
- aplicarea / adaptarea algoritmului Bayes Naiv pentru clasificare de texte:<sup>11</sup>  
ex. 6, ex. 26;  
folosirea regulii “add-one” [a lui Laplace] pentru „netezirea” parametrilor:  
ex. 6, ex. 27;

---

<sup>7</sup>De la cap. *Estimarea parametrilor; metode de regresie*, pentru estimarea parametrului unei distribuții Bernoulli vedeți ex. 1 și ex. 29.a, pentru estimarea parametrilor unei distribuții categoriale vedeți ex. 31, iar pentru estimarea parametrilor unei distribuții gaussiene vedeți ex. 7, ex. 38, ex. 8, ex. 39 (pentru cazul uni-variat) și ex. 10 (pentru cazul multi-variat).

<sup>8</sup>Vedeți cartea ML, pag. 156-157.

<sup>9</sup>Vedeți draftul capitolului suplimentar pentru cartea ML a lui T. Mitchell, *Generative and discriminative classifiers: Naive Bayes and logistic regression* (în special secțiunea 3).

<sup>10</sup>La secțiunea aceasta, precum și la următoarea secțiune, considerăm (implicit) că toate variabilele de intrare sunt de tip Bernoulli sau, mai general, de tip categorial. După aceea vom considera și variabile de intrare de tip continuu, în genere de tip gaussian. Variabila de ieșire se consideră întotdeauna de tip Bernoulli / categorial.

<sup>11</sup>Atenție: Noi am folosit aici versiunea de bază a algoritmului Bayes Naiv; varianta “bag of words” (vedeți cartea Machine Learning a lui Tom Mitchell, pag. 183) diferă ușor de aceasta.

- calculul ratei medii a erorilor pentru algoritmi Bayes Naiv și Bayes Corelat: ex. 10, ex. 11, ex. 28, ex. 29, ex. 30, ex. 31, ex. 35;
- evidențierea grafică a neconcordanței predicțiilor făcute de clasificatorii Bayes Naiv și Bayes Corelat: ex. 12.

## Proprietăți ale algoritmilor Bayes Naiv și Bayes Corelat

- (P0) 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;
- (P1) numărul de parametri necesari de estimat din date: liniar pentru Bayes Naiv ( $2d + 1$ ) și exponențial pentru Bayes Corelat ( $2^{d+1} - 1$ ): ex. 7.e, ex. 25.ab, ex. 30;
- (P2) 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];
- (P3) algoritmul Bayes Corelat poate produce eroare [la clasificare] din cauza faptului că ia decizia în sensul unui vot majoritar. Algoritmul Bayes Naiv are și el această „sursă” de eroare; în plus el poate produce eroare și din cauza faptului că lucrează cu presupuziția de independență condițională (care nu este satisfăcută în mod neapărat);
- (P4) acuratețea [la clasificare a] algoritmului Bayes Naiv scade atunci când unul sau mai multe atribute de intrare sunt duplicate: ex. 10.d, ex. 28.def;
- (P5) în cazul „învățării” unei funcții booleene (oarecare), rata medie a erorii produse la antrenare de către algoritmul Bayes Corelat (spre deosebire de Bayes Naiv!) este 0: ex. 30.d;
- (P6) complexitatea de eșantionare: de ordin logaritm pentru Bayes Naiv și de ordin exponențial pentru Bayes Corelat: ex. 13;
- (P7) corespondența dintre regula de decizie a algoritmului Bayes Naiv (când toate variabilele de intrare sunt de tip Bernoulli) și regula de decizie a *regresiei logistice* și, în consecință, liniaritatea granițelor de decizie: ex. 14.
- *comparații* între algoritmul Bayes Naiv și alți algoritmi de clasificare automată: ex. 33, ex. 35.

## Algoritmi Bayes Naiv și Bayes Corelat cu variabile de intrare de tip gaussian

- Aplicație: G[N]B: ex. 15, ex. 36, ex. 40;<sup>12</sup> GJB: ex. 37; GNB vs GJB: ex. 16.
- Proprietăți:

---

<sup>12</sup>Vedeți și ex. 38 de la cap. *Estimarea parametrilor; metode de regresie*.

- (P0') 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: ex. 17, ex. 36.a;
- (P1') similar, presupunem că variabila de ieșire este booleană; dacă variabilele de intrare (notație:  $X = (X_1, \dots, X_d)$ ) au distribuțiile [corelate] condiționale  $X|Y = 0$  și  $X|Y = 1$  de tip gaussian [multi-variabil], cu matricele de covarianță egale ( $\Sigma_0 = \Sigma_1$ ), atunci regula de decizie a algoritmului "full" / Joint Gaussian Bayes este și ea echivalentă (ca formă) cu cea a regresiei logistice, deci separarea realizată este tot de formă liniară: ex. 18;
- (P2') când variabilele de intrare satisfac condiții mixte de tip (P0') sau (P7), atunci concluzia – separare liniară – se menține: ex. 39.b;
- (P3') dacă în condițiile de la propozițiile (P0')-(P2') presupuziția de independență condițională este satisfăcută, iar numărul de instanțe de antrenament tinde la infinit, atunci rezultatul de clasificare obținut de către algoritmul Bayes Naiv gaussian este identic cu cel al regresiei logistice: ex. 19.a.  
Atunci când presupuziția de independență condițională nu este satisfăcută, iar numărul de instanțe de antrenament tinde la infinit, regresia logistică se comportă mai bine decât algoritmul Bayes Naiv [gaussian]: ex. 19.b;
- (P4') nu există o corespondență 1-la-1 între parametrii calculați de regresia logistică și între parametrii calculați de algoritmul Bayes Naiv [gaussian]: ex. 20.a;
- (P5') atunci când varianțele distribuțiilor gaussiene care corespund probabilităților condiționale  $P(X_i|Y = k)$  depind și de eticheta  $k$ , separatorul decizional determinat de algoritmul Bayes Naiv gaussian nu mai are forma regresiei logistice: ex. 38 (similar, pentru algoritmul Bayes Corelat gaussian, atunci când  $\Sigma_0 \neq \Sigma_1$ : ex 37);
- (P6') parametrii algoritmului Bayes Corelat gaussian se pot estima în timp liniar în raport cu numărul de instanțe din setul de date de antrenament: ex. 20.b.



## 5. Învățare bazată pe memorare

### Sumar

#### Noțiuni preliminare

- măsuri de distanță, măsuri de similaritate: ex. 2;
- normă într-un spațiu vectorial; [măsura de] distanță indusă de către o normă: ex. 7;
- $k$ -NN vecinătate a unui punct din  $\mathbb{R}^d$ .

#### Algoritmul $k$ -NN

- pseudo-cod: cartea ML, pag. 232;
- bias-ul inductiv: „Cine se aseamănă se adună” (sau: „Spune-mi cu cine te împrietenești, ca să-ți spun cine ești”): ex. 15.a;
- exemple (simple) de aplicare: ex. 1, ex. 2, ex. 15.b-d;
- complexitate de spațiu:  $\mathcal{O}(dn)$   
complexitate de timp:
  - la antrenare:  $\mathcal{O}(dn)$
  - la testare:  $\mathcal{O}(dn \log n)$

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

unde  $d$  este numărul de attribute, iar  $n$  este numărul de exemple;
- arbori  $kd$  (engl.,  $kd$ -trees): *Statistical Pattern Recognition*, Andrew R. Webb, 3rd ed., 2011, Willey, pag. 163-173;
- $k$ -NN ca algoritm ML “lazy” (vs. “eager”):  
suprafețe de decizie și granițe de decizie:  
diagrame Voronoi pentru 1-NN: ex. 4, ex. 11.a, ex. 18, ex. 19, ex. 20.a;
- analiza erorilor:
  - 1-NN pe date consistente: eroarea la antrenare este 0: ex. 2, ex. 12.a;
  - variația numărului de erori (la antrenare și respectiv testare) în funcție de valorile lui  $k$ : ex. 22, ex. 23.ab;
  - $k$ -NN ca metodă ne-parametrică; alegerea lui  $k$ : CV: ex. 23.c;
  - CVLOO: ex. 3, ex. 12.b, ex. 16, ex. 24.a, ex. 20.b;
  - sensibilitatea / robustețea la „zgomote”: ex. 5, ex. 15;
  - eroarea asimptotică: ex. 10, ex. 25.
- efectul trăsăturilor redundante sau irelevante;
- alegerea valorii convenabile pentru  $k$ : ex. 21.

#### Proprietăți ale algoritmului $k$ -NN

- (P0) output-ul algoritmului  $k$ -NN pentru o instanță oarecare de test  $x_q$  depinde de valoarea lui  $k$ : ex. 1;

- (P1) pe seturi de date de antrenament *consistente*, eroarea la antrenare produsă de algoritmul 1-NN este 0: ex. 2, ex. 12.a;
- (P2) output-ul algoritmului  $k$ -NN, precum și suprafețele de decizie și separatorii decizionali depind de *măsura de distanță* folosită: ex. 7;
- (P3) „blestemul marilor dimensiuni” (engl., the curse of dimensionality): în anumite condiții, numărul de instanțe de antrenament necesare pentru a avea un *cel mai apropiat vecin* situat la distanță *rezonabilă* față de instanța de test  $x_q$  crește exponențial în funcție de numărul de attribute folosite: ex. 9;
- (P4) în anumite condiții, rata medie a *erorii asimptotice* a algoritmului 1-NN este mărginită superior de dublul ratei medii a erorii algoritmului Bayes Corelat: ex. 10, ex. 25.

### Comparații cu alți algoritmi de clasificare automată

- ID3: ex.11.b, ex. 13.ab;
- SVM: ex.12, ex. 13.c, ex. 24.b;
- regresia logistică: ex. 24.b;
- 1-NN cu mapare cu RBF: ex. 14.

### Variante ale algoritmului $k$ -NN

- $k$ -NN folosind alte măsuri de distanță (decât distanța euclidiană): ex. 7;
- $k$ -NN cu *ponderarea distanțelor* (engl., distance-weighted  $k$ -NN): cartea ML, pag. 236-238 (formulele 8.2, 8.3, 8.4);<sup>13</sup>
- algoritmul lui Shepard: ex.8.

### Alte metode de tip IBL

- rețele RBF: cartea ML, pag. 238-240;
- raționare bazată pe cazuri (engl., case-based reasoning): cartea ML, pag. 240-244.

---

<sup>13</sup>Sectiunea 8.3 din cartea ML (pag. 236-238) se referă la regresia [liniară] local-ponderată ca o formă mai generală de aproximare a [valorilor] funcțiilor, în raport cu cele calculate de către algoritmul  $k$ -NN atunci când se folosește ponderarea distanțelor.

## 6. Clusterizare

### Sumar

#### 6.0 Noțiuni de bază

- instanță neetichetată vs. instanță etichetată (exemplu de antrenament);
- învățare nesupervizată (clusterizare) vs. învățare supervizată (clasificare);
- [funcție / măsură de] *distanță* definită pe  $\mathbb{R}^d \times \mathbb{R}^d$ : ex. 2 de la capitolul *Învățare bazată pe memorare*;
- cluster / grup / grupare / bin (engl.) vs. clasă;
- tipuri de clusterizare: ierarhică vs. neierarhică;
- tipuri de ierarhii: ierarhii (arbori de clusterizare, dendrograme) obișnuite vs. ierarhii plate (engl., flat hierarchies);  
exemple: ex. 1.a și respectiv ex. 1.b, ex. 6.a;
- tipuri de apartenență a unei instanțe la un cluster: hard vs. soft (ultima numai pt. clusterizare ne-ierarhică).

#### 6.1. Clusterizare ierarhică

##### 6.1.1. 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”:<sup>14</sup>  $d(A, B) = \min\{d(x, y) | x \in A, y \in B\}$

“complete-linkage”:<sup>15</sup>  $d(A, B) = \max\{d(x, y) | x \in A, y \in B\}$

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

*metrica lui Ward*: ex. 30, 31.

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 [al algoritmului de clusterizare ierarhică] și orice alt cluster  $C$ ;

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

exemplu (pentru clustere non-singleton):

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

---

<sup>14</sup>Sau: nearest-neighbour.

<sup>15</sup>Sau: furthest-neighbour.

### 6.1.2. 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: ex. 3.b;
- exemple de aplicare: ex. 1-5, ex. 25-29 (pentru bottom-up), respectiv ex. 6 (pentru top-down);
- implementări: ex. 33, ex. 31, ex. 34.

### 6.1.3 Proprietăți

- (P0) 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: ex. 5 și ex. 28;
- (P1) 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: ex. 4.a; numărul minim de niveluri:  $\lceil \log_2 n \rceil$ ; ex. 4.b;
- (P2) există o anumită corespondență între clusterizare ierarhică cu similaritate de tip
  - “single-linkage” și aflarea *arborelui* [de acoperire] de cost minim dintr-un graf: ex. 6;
  - “complete-linkage” și aflarea unei *clici* (subgraf maximal complet) dintr-un graf (vedeți Manning & Schütze, *op. cit.*, pag. 506-507);
- (P3) algoritmul de clusterizare aglomerativă la al cărui pseudo-cod am făcut referire mai sus are complexitate  $\mathcal{O}(n^3)$ : ex. 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”:  
(P4) dacă se folosește ca măsură de similaritate între 2 instanțe cosinusul unghiului dintre vectorii care reprezintă 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] unui cluster nou format, precum și calculul „distanței” dintre două clustere se pot face în timp constant: ex. 32.

## 6.2. Clusterizare neierarhică,

*folosind asignare “hard” a instanțelor la clustere*

### 6.2.1 Noțiuni specifice

- centroid (centru de greutate) al unui cluster,  
 $K$ -partiție,  $K$ -configurație [inițială] a centroizilor: ex. 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,  $\mu$  este  $K$ -configurație de centroizi, iar  $\mu_{C(x_i)}$  este centroidul cel mai apropiat de  $x_i$ : ex. 12.

### 6.2.2 Algoritmul $K$ -means

- pseudo-cod (o versiune [mai] generală): Manning & Schütze, *op. cit.*, pag. 516; alternativ, vedeți enunțul ex. 12 (sau, echivalent, folosind variabile-indicator: ex. 39);  
exemple de aplicare: ex. 7-11, ex. 15.a, ex.19.a, ex. 20.a, ex. 35, ex. 36.
- exemple de *euristici pentru inițializarea centrozilor*:  
inițializare arbitrară / random în  $\mathbb{R}^d$  sau în  $X = \{x_1, x_2, \dots, x_n\} \subseteq \mathbb{R}^d$  (setul de date de clusterizat);  
aplicare în prealabil a unui algoritm de clusterizare ierarhică;  
folosind o anumită distribuție probabilistă definită pe  $X$ :  $K$ -means++ (David Arthur, Sergei Vassilvitskii, 2007): ex. 43.
- 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 centrozilor 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  $\varepsilon$  fixat în prealabil.
- ca algoritm *per se*:  
 $K$ -means este un algoritm de *căutare*:  
*spațiul de căutare* este mulțimea tuturor  $K$ -partițiilor care se pot forma pe dataset-ul de intrare;  
(P0) întrucât acest spațiu de căutare (deși este finit) este exponențial ( $K^n$ ),  $K$ -means explorează doar parțial spațiul de căutare, procedând *iterativ*: el pleacă de la o „soluție” ( $K$ -partiție) aleasă eventual în mod arbitrar / aleatoriu și o „îmbunătățește” la fiecare iterație;  
(P1) soluția găsită este dependentă de inițializarea centrozilor: ex. 10;  
(P1') 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: ex. 12.b;  
(P1'') rezultatele lui  $K$ -means sunt dependente [și] de măsura de distanță folosită: ex. 42.  
 $K$ -means poate fi văzut și ca *algoritm de optimizare* — vedeți criteriul  $J_K$  de mai sus;  
(P2) 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  $J_K(C^t, \mu^t)$ : ex. 12.a;  
(P2') algoritmul  $K$ -means nu garantează atingerea optimului global (i.e., minimul) criteriului  $J_K$ : ex. 12.b, ex. 40.b.
- ca algoritm de *învățare automată*:  
[urmat de] „generalizare”: o instanță nouă  $x$  se asociază clusterului având

centroidul cel mai apropiat de  $x$ ;

(P3) „granițele“ de separare dintre [perechile de] clustere produse de  $K$ -means sunt [doar] liniare, [cel puțin] atunci când se folosește distanța euclidiană: ex. 11.b;

(P3') este însă posibil să se obțină separatori neliniari dacă se folosește o versiune „kernelizată“ a algoritmului  $K$ -means: ex. 44;

(P4) rezultatele lui  $K$ -means pot fi influențate de prezența outlier-elor: ex. 10.

- chestiunea alegerii unei valori convenabile / „naturale” pentru  $K$  (pentru un dataset dat): ex. 38 (și CMU, 2012f, E. Xing, A. Singh, HW3, ex. 1.de).
- adaptarea algoritmului  $K$ -means pentru cazul în care în locul distanței euclidiene se folosește distanța Manhattan: ex. 42;
- implementare: ex. 47.

### 6.2.3 Alte proprietăți ale algoritmului $K$ -means

- 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,  $X = \{x_1, x_2, \dots, x_n\} \subseteq \mathbb{R}^d$ :
  - (P5) pentru  $K > 0$  fixat,  $|\mathcal{P}_K| = K^n$ , deci este finit, și există  $\underline{J}_K \stackrel{not.}{=} \min_C J_K(C, \mu_C)$ ;
  - acest minimum ( $\underline{J}_K$ ) se poate obține prin explorarea exhaustivă a spațiului  $\mathcal{P}_K$ , însă consumul de timp este prohibitiv în practică: ex. 12.b;
  - (P6) valoarea 0 pentru  $\underline{J}$  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  $i = 1, \dots, n$  (ex. 38);
  - (P7)  $\underline{J}_1 \geq \underline{J}_2 \geq \dots \geq \underline{J}_{n-1} \geq \underline{J}_n = 0$ : ex. 13.
- (P8) 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”:  $C_1 = \{x_1, \dots, x_{i_1}\}$ ,  $C_2 = \{x_{i_1+1}, \dots, x_{i_2}\}$ ,  $\dots$ ,  $C_K = \{x_{i_{K-1}+1}, \dots, x_n\}$ , cu  $i_1 < i_2 < \dots < i_{K-1} < i_K = n$ ;
  - există un algoritm [de programare dinamică] de complexitate  $\mathcal{O}(Kn^2)$  care calculează  $\underline{J}_K$ : ex. 40.
- în legătură cu  $J_K$  și algoritmul  $K$ -means:
  - (P9)  $J_K(C^{t-1}, \mu^{t-1}) \geq J_K(C^t, \mu^t)$  la orice iterație ( $t > 0$ ) a algoritmului  $K$ -means: ex. 12.a;
  - (P9') 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;
  - (P10) în vreme ce minimizează *coeziunea intra-clustere*, i.e. o variantă ponderată a „sumelor celor mai mici pătrate” calculate pe clustere,

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

unde  $\gamma_{ik} = 1$  dacă  $x_i$  aparține clusterului de centroid  $\mu_k$  și  $\gamma_{ik} = 0$  în caz contrar, algoritmul  $K$ -means maximizează (în mod aproximativ!) o sumă ponderată a distanțelor dintre 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  $x_1, x_2, \dots, x_n$  (ex. 39).

### 6.3. Clusterizare neierarhică, folosind asignare “soft” a instanțelor la clustere

#### 6.3.1 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: ex. 18 de la capitolul de *Fundamente*;
- distribuție (funcție de densitate) de probabilitate (p.d.f.);  
parametri ai unei distribuții de probabilitate;  
distribuția gaussiană: cazurile uni- și multi-variat;
- mixtură de distribuții probabiliste:  
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, respectiv o mixtură de distribuții Bernoulli: ex. 23 și ex. 85 de la capitolul de *Fundamente*; o mixtură de distribuții gaussiene multi-variate: ex. 89 de la capitolul de *Fundamente*; o mixtură de distribuții oarecare: ex. 90 de la capitolul de *Fundamente*;
- funcție de *verosimilitate* a unui set de date ( $D$ ), în raport cu o distribuție probabilistă dată:  $L(\theta) = P(D|\theta)$ , unde prin  $\theta$  se notează parametrii respectivei distribuții. Exemplificare: ex. 1.abd, ex. 2 de la capitolul *Estimarea parametrilor; metode de regresie*;
- MLE (Maximum Likelihood Estimation): estimarea [valorilor] parametrilor unei distribuții probabiliste în sensul maximizării verosimilității datelor disponibile. Exemplificare: capitolul *Estimarea parametrilor; metode de regresie*, ex. 1-11, ex. 29-40. Aplicare în cazul distribuției gaussiene uni-variate: ex. 15.ab de la capitolul *Clasificare bayesiană*;
- analiza discriminativă gaussiană: ex. 38 de la capitolul *Estimarea parametrilor; metode de regresie*;
- Observație: Algoritmul EM este [sau, mai degrabă, poate fi folosit ca] o metodă de estimare a parametrilor unei mixturi de distribuții probabiliste. *Alternativ*, pentru același obiectiv pot fi folosite alte metode, de exemplu *metoda gradientului ascendent*: ex. 59.

### 6.3.2 Algoritmul EM pentru clusterizare prin estimarea parametrilor unui model de mixturi de distribuții gaussiene (EM/GMM)

- pseudo-cod:
  - cazul uni-dimensional, varianta când doar parametrul  $\mu$  este lăsat liber: ex. 48 (cf. *Machine Learning*, Tom Mitchell, 1997, pag. 193); aplicare: ex. 15.b, ex. 16;
  - cazul uni-dimensional, varianta când toți parametrii ( $\pi$ ,  $\mu$  și  $\sigma$ ) sunt lăsați liberi: ex. 17 (aplicare: ex. 15.c);
  - alte variante: ex. 49, ex. 50, ex. 51;
  - cazul multi-dimensional, varianta când toți parametrii ( $\pi$ ,  $\mu$  și  $\Sigma$ ) sunt lăsați liberi: ex. 23; alte variante: ex. 52, ex. 54;
  - aplicarea algoritmului EM/GMM, cazul bivariat: ex. 19.b, ex. 20.b, ex. 21, ex. 22, ex. 55, ex. 56, ex. 57;
- schema algoritmică EM: vedeți Tom Mitchell, *Machine Learning* book, 1997, pag. 194-195;
- ca algoritm de *învățare statistică*:
  - algoritmul EM poate fi văzut ca o metodă de estimare a parametrilor (engl., parameter fitting);
- ca algoritm *per se*:
  - *algorithm 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;
  - *algorithm de optimizare*:

în *esență* / *rezumat*, metoda de maximizare a funcției de *log-verosimilitate* a datelor observabile  $\log P(X|\theta)$  este maximizarea la fiecare iterație  $t$  a unei funcții auxiliare  $Q_t$ , care constituie o margine inferioară a lui  $\log P(X|\theta)$ , și anume media funcției de *log-verosimilitate* a datelor complete în raport cu distribuția de probabilitate a variabilelor neobservabile la iterația  $t$ ;

așadar, la fiecare iterație  $t$  se calculează funcția „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 ( $\theta$ ), se alege în mod arbitrar, iar apoi  $\theta^{(t+1)} = \operatorname{argmax}_{\theta} Q_t(\theta|\theta^{(t)})$ ;

media reprezentată de funcția  $Q_t$  se calculează în funcție de distribuțiile condiționale ale variabilelor „neobservabile“  $Z$  în raport cu datele observabile  $X$  și cu  $\theta^{(t)}$ ;

(P0) Se poate demonstra că funcția  $Q_t$  constituie o *margine inferioară* pentru funcția de log-verosimilitate a variabilelor „observabile“,  $\log P(X|\theta)$ : ex. 1 de la capitolul *Algoritmul EM*;

(P1) *Teorema de corectitudine* (vedeți ex. 1 și în special ex. 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;



- ca algoritm de *învățare automată*:  
algoritmul EM este o metodă de identificare / învățare de ipoteze ML (Maximum Likelihood); vedeți 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  $\max_{j'} P(X = x|h_{j'})P(h_{j'})$ ;  
(P2) Rezultatele algoritmului EM depind (ca și la  $K$ -means) de valorile atribuite parametrilor la inițializare (ex. 15.c).  
(P3) Anumite valori atribuite inițial parametrilor algoritmului EM pot provoca rularea la infinit a algoritmului, fără ca [la pasul M] valorile parametrilor să se modifice de la o iterație la alta: ex. 18.c;  
(P4) Spre deosebire de cazul algoritmului  $K$ -means, suprafețele / granițele de separare create de algoritmul EM/GMM nu sunt în mod neapărat liniare (vedeți de exemplu situațiile întâlnite la rezolvarea ex. 15.c, pag. 556, sau la ex. 57.c și ex. 58.c).
- (P5) Comparativ cu algoritmul  $K$ -means, algoritmul EM/GMM este în general mai lent — mișcarea centroizilor poate explora într-o manieră mai fină spațiul (vedeți de exemplu ex. 19) —, iar din acest motiv el poate să obțină uneori rezultate mai bune / convenabile (vedeți spre exemplu ex. 20); EM/GMM este mai robust la influența outlier-elor;  
(P6) Apare un fenomen de “atracție” reciprocă a mediilor gaussianelor (aceste medii fiind echivalentul centroizilor din algoritmul  $K$ -means), datorită faptului că fiecare instanță aparține (cu o anumită probabilitate) la fiecare cluster. Atracția mediilor este cu atât mai puternică cu cât varianțele sunt mai mari. (Vedeți spre exemplu ex. 15.b.)

### 6.3.3 Alte proprietăți ale algoritmului EM/GMM

- Pentru distribuții gaussiene multi-variate:
  - (P7) în cazul cel mai general (deci când matricea  $\Sigma$  nu este neapărat diagonală), datele generate de acest tip de distribuție se grupează în elipse (corpuri elipsoidale) cu axe de simetrie [desigur, perpendiculare, dar altfel] nerestricționate.
  - (P7') dacă matricea de covarianță  $\Sigma$  este diagonală, atunci distribuția gaussiană respectivă este echivalentă cu un set / vector de variabile gaussiene uni-variate independente: ex. 28 de la capitolul de *Fundamente*;
  - (P7'') dacă matricea  $\Sigma$  este de forma  $\sigma^2 I$ , unde  $I$  este matricea identitate, datele generate de respectiva distribuție tind să se grupeze în sfere;
  - (P7''') dacă matricea  $\Sigma$  este diagonală (fără nicio altă restricție), datele generate se grupează în elipse (sau: corpuri elipsoidale) având axe de simetrie paralele cu axele sistemului de coordonate;
- (P8) 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: ex. 54;
- O legătură interesantă între algoritmul EM/GMM și metoda gradientului ascendent, în cazul în care matricele de covarianță sunt de forma  $\sigma_k^2 I$ : ex. 59;

- O legătură interesantă între clasificatorul Bayes Naiv gaussian și algoritmul EM/GMM în cazul în care matricele de covarianță sunt de forma  $\sigma_k^2 I$ : o variantă semi-supervizată a algoritmului EM/GMM: ex. 60.
- Schema algoritmică EM (vedeți Tom Mitchell, *Machine Learning* book, 1997, pag. 194-195) are diverse variante și aplicații:
  - calculul parametrilor pentru mixturi de diverse distribuții [nu doar gaussiene]: vedeți 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 schema algoritmică EM: vedeți cele menționate mai sus în legătură cu algoritmul EM văzut ca *algoritm de optimizare*.

## 7. Algoritmul EM

### Sumar

#### Noțiuni preliminare

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

#### Schema algoritmică EM

- pseudo-cod: *Machine Learning*, Tom Mitchell, 1997, pag. 194-195;
- fundamentare teoretică: ex. 1 — pentru funcția  $F(q, \theta)$ , care reprezintă o margine inferioară a funcției de log-verosimilitate a datelor complete, se face maximizarea prin metoda creșterii pe coordonate (engl., coordinate ascent) — și ex. 2 (monotonia funcției de log-verosimilitate a datelor complete);<sup>16</sup>
- chestiuni metodologice (relativ la inițializarea parametrilor): ex. 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*: ex. 14 și ex. 4;
  - mixturi de distribuții *catoriale*: ex. 5 și ex. 15;
  - mixturi de distribuții *Poisson*: ex. 18;
  - mixturi de distribuții *Gamma*: ex. 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”: ex. 3;  
similar pentru distribuția *multinomială*: ex. 13;
- EM pentru estimarea tuturor parametrilor unei distribuții *categoriale*: ex. 12;
- EM pentru estimarea parametrului unei distribuții *Poisson* în condițiile în care o parte din valorile date lipsesc: ex. 10;

---

<sup>16</sup> Legătura cu funcția auxiliară de la iterația  $t$ :  $Q(\theta \mid \theta^{(t)}) = \underbrace{F(P(z \mid x, \theta^{(t)}), \theta)}_{G_t(\theta)} - H[P(z \mid x, \theta^{(t)})]$ .

- EM pentru estimarea parametrilor a două distribuții probabiliste atunci când se dau instanțe care sunt generate de *suma* celor două distribuții:  
distribuții *exponențiale*: ex. 8,  
distribuții *gaussiene*: ex. 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): ex. 7 (varianta de asignare “soft” a instanțelor la cluster) și ex. 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: ex. 9;
- EM pentru *modelul mixturii domeniilor semantice* (engl., topic model) pentru clusterizare de *documente*: [ex. 6 și] ex. 16;
- EM pentru *estimarea* [nu în sens MLE, cum a fost cazul până aici, ci] *în sens MAP*: ex. 20.

## 8. Rețele neuronale artificiale

### Sumar

#### 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;
- metoda gradientului descendent (ca metoda de optimizare); avantaje și dezavantaje; ex. 54 de la cap. *Fundamente*, ex. 22, ex. 34, ex. 35.

#### 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 conexiunilor 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).

#### Câteva proprietăți relative la *expresivitatea* 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: ex. 8.

- (P1) Rețele de neuroni diferite (ca structură și / sau tipuri de unități) pot să calculeze o aceeași funcție: ex. 3 și ex. 1.c vs. ex. 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: ex. 4 vs. ex. 33.<sup>17</sup>
- (P2) Orice unitate liniară situată pe un nivel ascuns poate fi „absorbită” pe nivelul următor: ex. 32.
- (P3) Orice funcție booleană poate fi reprezentată cu ajutorul unei rețele neuronale artificiale având doar două niveluri de perceptroni-prag: ex. 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: ex. 7.

### Algoritmi de antrenare a neuronilor artificiali folosind metoda gradientului descendent

- algoritmul de antrenare a unității liniare: ex. 36;  
vedeți T. Mitchell, *Machine Learning*, p. 93, justificare: p. 91-92; convergența: p. 95; exemplu de aplicare: ex. 10;  
variante incrementală a algoritmului de antrenare a unității liniare: cartea ML, p. 93-94; despre convergența acestei variante (ca aproximare a variantei precedente (“batch”)): cartea ML, p. 93 jos;
- algoritmul de antrenare a perceptronului-prag și convergența: cartea ML, p. 88-89; exemplu de aplicare: ex. 11;
- algoritmul de antrenare a perceptronului sigmoidal și justificarea sa teoretică: cartea ML, p. 95-97;
- algoritmul *Perceptron* al lui Rosenblatt; exemplu de aplicare: ex. 16, ex. 38;
- deducerea regulii de actualizare a ponderilor pentru tipuri particulare de perceptroni: ex. 12, ex. 24.a, ex. 37, ex. 13.a;
- o justificare probabilistă (gen ipoteză de tip *maximum likelihood*) pentru minimizarea sumei pătratelor erorilor [la deducerea regulii de antrenare] pentru perceptronul liniar: ex. 13.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, ex. 14 (pentru perceptronul liniar), o funcție de tip cross-entropie, ex. 15 (pentru perceptronul sigmoidal).

### Perceptronul Rosenblatt și rezultate de *convergență*

- exemplu de aplicare [adică, învățare cu perceptronul Rosenblatt]: ex. 16.
- câteva *proprietăți* simple ale perceptronului Rosenblatt: ex. 17.

<sup>17</sup>Problemele 1.d și ex. 31 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.

- rezultate de convergență de tip “mistake bound” pentru [algoritmul de antrenare pentru] perceptronul-prag [în varianta] Rosenblatt: ex. 18, ex. 39; pentru perceptronul-prag (clasic): ex. 41; învățare online cu perceptronul-prag de tip Rosenblatt: ex. 40;
- *Perceptronul* kernel-izat [dual]: ex. 23; particularizare pentru cazul nucleului RBF: ex. 49.

### Antrenarea rețelelor neuronale artificiale:

#### algoritmul de *retro-propagare* pentru rețele 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 al rețelelor feed-forward (de unități sigmoidale) cu oricâte niveluri, vedeți p. 101-103; ex. 19: deducerea regulilor de actualizare a ponderilor în cazul rețelelor cu 2 niveluri, având însă unități cu funcție de activare oarecare (derivabilă);
- aplicare: ex. 20, ex. 42, ex. 43;
- prevenirea overfitting-ului: folosirea unei componente de tip „moment” în expresia regulilor de actualizare a ponderilor: ex. 45; regularizare: introducerea unei componente suplimentare în funcția de optimizat: ex. 21;
- cazul folosirii unei funcții de activare de tip tangentă hiperbolică: ex. 44;
- cazul folosirii unei funcții de cost / penalizare / eroare de tip cross-entropie: ex. 47;
- 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: ex. 48.

### Rețele neuronale profunde — câteva chestiuni introductive

- analiza convexității unor funcții de cost folosite în învățarea profundă: ex. 54;
- 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: ex. 26;
- determinarea numărului de parametri și de conexiuni din rețeaua neuronală convolutivă LeNet: ex. 27;
- determinarea mărimii hărții de trăsături de pe un anumit nivel, precum și a numărului de operații în virgulă mobilă (FLOPs) executate la procesarea forward într-o rețea neuronală convolutivă: ex. 55.

## 9. Mașini cu vectori-suport

### Sumar

#### Noțiuni preliminare

- elemente [simple] de *calcul vectorial*; proprietăți elementare ale produsului scalar al vectorilor din  $\mathbb{R}^n$ , norma euclidiană ( $L_2$ ) și norma  $L_1$  în  $\mathbb{R}^n$ : ex. 1, ex. 34;
- elemente [simple] de *geometrie analitică*: ecuația unei drepte din planul euclidian, ecuația unui plan din  $\mathbb{R}^3$ , ecuația unui hiper-plan din  $\mathbb{R}^n$ ; ecuația dreptei care trece prin două puncte date în planul euclidian: ex. 5.c; panta unei drepte perpendiculare pe o dreaptă dată: ex. 9.d, ex. 35, ex. 37;
- distanța (cu sau fără semn) de la un punct la o dreaptă (respectiv la un plan, sau mai general la un hiperplan): ex. 5, ex. 6.c, ex. 1;
- proprietăți de bază din *calculul matriceal*;
- calculul *derivatelor parțiale* [pentru funcții obținute prin compuneri de funcții elementare];
- *metoda lui Lagrange* pentru rezolvarea problemelor de *optimizare convexă cu restricții*: ex. 34, ex. din [https://www.cs.helsinki.fi/u/jkivinen/teaching/sumale/Spring2014/kkt\\_example.pdf](https://www.cs.helsinki.fi/u/jkivinen/teaching/sumale/Spring2014/kkt_example.pdf), ex. 56, ex. 57 de la capitolul de *Fundamente*;
- *separabilitate* liniară, separator optimal, *margină* geometrică: ex. 37, ex. 6.abc, ex. 9.

#### SVM cu margine “hard”

- (•) deducerea *forme primale* pentru problema SVM (cu margine “hard”), pornind de la principiul maximizării marginii geometrice: ex. 2;<sup>18</sup>
- (P0) o formă [simplă] echivalentă cu *forma primală a problemei de optimizare SVM*: ex. 7;
- *exemplificarea* identificării *separatorului optimal* și a *vectorilor-suport*, pornind de la condițiile din forma primală: ex. 3-5, ex. 35, ex. 37, ex. 38 și CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, HW4, ex. 4.1-4;
- calcularea *erorii* la cross-validare “leave-one-out” atunci când se folosește o SVM liniară cu margine “hard”: ex. 4.d, ex. 36;
- exemple de [funcții de] *mapare a atributelor*, cu scopul de a obține separabilitate liniară: ex. 6.d, ex. 8, ex. 9.b, ex. 9, ex. 39.bd, ex. 40.a; rezolvarea directă a problemei SVM primale în [noul] spațiu de trăsături; identificarea separatorului neliniar din spațiul inițial [de trăsături]: ex. 9.de, ex. 39.ce, ex. 40.b-e;
- (•) deducerea *forme duale* pentru problema SVM cu margine “hard”: ex. 10;

---

<sup>18</sup>Vedeți și Andrew Ng (Stanford), Lecture Notes, part V, section 3.



- *exemplificarea* a două modalități de găsim a formei duale a problemei SVM cu margine “hard” pentru învățarea unui concept [reprezentat de un set de date de antrenament neseparabil în spațiul „inițial” de trăsături] folosind o funcție de mapare  $\Phi$  dată: prin optimizare directă (ex. 11, ex. 41), respectiv prin folosirea relațiilor de legătură cu soluția problemei primale: ex. 12;
- (P1) efectul *multiplicării valorilor atributelor* cu o constantă pozitivă asupra separatorului obținut de SVM): ex. 27.a;<sup>19</sup>
- vedeți proprietățile (P4) și (P6) enunțate mai jos (la secțiunea despre C-SVM), care sunt valabile și în cazul SVM [cu margine “hard”];
- (P2) efectul unui *atribut irelevant* — în sensul că nu afectează satisfacerea *restricțiilor* de separabilitate liniară a datelor de antrenament și, în plus, nu mărește marginea de separare — asupra rezultatelor clasificatorului SVM (și respectiv C-SVM): ex. 48, ex. 56;
- *comparații* între SVM [având, eventual, diferite funcții-nucleu] și alți clasificatori: ex. 56, ex. 44, ex. 45; vedeți și ex. 12 de la capitolul *Învățare bazată pe memorare*.

### SVM cu margine “soft” (C-SVM):

- (•) deducerea *formei duale* pentru problema SVM cu margine “soft” (C-SVM): ex. 13;
- (P3) exprimarea / calcularea distanței geometrice de la un vectori-suport  $x_i$  pentru care  $\bar{\alpha}_i = C$  la hiperplanul-margine corespunzător etichetei  $y_i$ , cu ajutorul variabilei de „destindere”  $\xi_i$ : ex. 14.a;
- *exemplificarea* noțiunilor de bază: ex. 46, ex. 47 și CMU, 2008 fall, Eric Xing, final, ex. 2.2;  
un *exemplu* de calculare a valorii optime pentru funcția obiectiv a problemei de optimizare C-SVM: ex. 14.b;
- exemplificarea poziționării separatorului optimal determinat de C-SVM (pentru diferite valori ale parametrului  $C$ ), în prezența unui outlier: ex. 16;
- exemplificarea efectului pe care îl are creșterea valorii parametrului de „destindere”  $C$  [asupra marginii și asupra excepțiilor la clasificare]: ex. 17, CMU, 2010 fall, Aarti Singh, HW3, ex. 3.2;
- un exemplu de situație în care forma duală a problemei de optimizare C-SVM are soluție unică, dar forma sa primală *nu* are soluție unică: ex. 18;
- (P4) o proprietate pentru C-SVM (dar și pentru SVM): ex. 15  
Dacă în setul de date de antrenament două trăsături (engl., features) sunt duplicate ( $x_{ij} = x_{ik}$  pentru  $i = 1, \dots, m$ ), atunci ele vor primi ponderi identice ( $\bar{w}_j = \bar{w}_k$ ) în soluția optimală calculată de clasificatorul [C-]SVM;
- (P5) o *margine superioară* (engl., upper bound) pentru numărul de *erori* comise la *antrenare* de către C-SVM: ex. 19;

$$\text{err}_{\text{train}}(\text{C-SVM}) \leq \frac{1}{m} \sum_i \xi_i$$

---

<sup>19</sup> Rezultatul *nu* se menține și în cazul C-SVM: ex. 27.b.

- (P6) o proprietate pentru C-SVM (dar și pentru SVM):  
La CVLOO numai vectorii-suport pot fi (eventual!) clasificați eronat: ex. 20;  
așadar avem [și] o *margină superioară* pentru numărul de *erori* comise la CVLOO de către C-SVM:

$$\text{err}_{CVLOO}(\text{C-SVM}) \leq \frac{\#SV_s}{m}$$

- (P7) chestiuni legate de *complexitatea computațională* privind clasificatorul C-SVM: ex. 54;
- (•) deducerea *formeî duale* pentru problema SVM cu margine “soft” (C-SVM) de normă  $\mathcal{L}_2$ : ex. 49;
- (•) o formă echivalentă a problemei de optimizare C-SVM, în care nu apar deloc restricții asupra variabilelor, dar în care se folosește funcția de pierdere / cost *hinge*: ex. 21;  
exemplificare / aplicare (și comparare cu regresia logistică): ex. 50;
- (•) algoritmul SMO (Sequential Minimal Optimization): deducerea relațiilor de actualizare a variabilelor Lagrange;  
*exemple* de aplicare a algoritmului SMO simplificat: ex. 23, ex. 51;
- o *comparație* asupra efectului *atributelor irelevante* (aici, în sensul că odată eliminate / adăugate, n-ar trebui să afecteze rezultatele clasificării) asupra clasificatorilor 1-NN și C-SVM: ex. 56;

## SVM / C-SVM și funcțiile-nucleu — câteva proprietăți

- *exemplificarea* corespondenței dintre forma (primală sau duală) a problemei C-SVM și alegerea valorii parametrului de „destindere”  $C$  și a *funcției-nucleu* pe de o parte și alura și poziționarea separatorului optimal pe de altă parte: ex. 24, ex. 52;
- *exemplificarea* efectului pe care îl are translatarea datelor în raport cu o axă ( $Oy$ ) asupra poziției separatorului optimal (în raport cu *funcția-nucleu* folosită): ex. 42;
- C-SVM: condiții suficiente asupra parametrului de „destindere”  $C$  și asupra valorilor funcției-nucleu pentru ca toate instanțele de antrenament să fie vectori-suport: ex. 25;
- SVM cu nucleu RBF: câteva proprietăți remarcabile
  - pentru SVM pe un set de date [separabil liniar în spațiul de trăsături] instanțe foarte depărtate de separatorul optimal pot fi vectori-suport: ex. 43;
  - (P8) pentru orice set de instanțe distincte și pentru orice etichetare a acestora, există o valoare a hiper-parametrului nucleului RBF ( $\sigma$ ) astfel încât SVM obține la antrenare eroare 0: ex. 26. Rezultatul *nu* este valabil și pentru C-SVM;
  - (P9) pentru orice set de instanțe distincte, pentru orice etichetare a acestora și pentru *orice* valoare a hiper-parametrului nucleului RBF ( $\sigma$ ), problema de tip SVM care impune ca toate instanțele să fie corect clasificate și la distanța  $1/\|w\|$  față de separatorul optimal are soluție: ex. 58;
- avantaje și dezavantaje ale folosirii metodelor de clasificare liniară de tipul SVM, Perceptron etc. și versiunile lor kernel-izate: ex. 55.

- chestiuni recapitulative: ex. 28, ex. 27, ex. 66, ex. 57.

### Alte probleme [de optimizare] de tip SVM

- SVM pentru clasificare  $n$ -ară (SVM multi-class): ex. 29, ex. 59;
- deducerea *formeii duale* pentru problema *one-class SVM*, versiunea *Max Margin*: ex. 30 (varianta cu margine “hard”), ex. 61 (varianta cu margine “soft”, folosind  $\nu$ -SVM);
- legătura dintre soluțiile problemei *one-class SVM*, versiunea *Max Margin*, cu margine “hard” și respectiv cele ale problemei SVM (cu și respectiv fără termen liber (engl., bias), tot cu margine “hard”: ex. 60;
- deducerea *formeii duale* pentru problema *one-class SVM*, versiunea *minimum enclosing ball* (MEB): ex. 31 (varianta cu margine “hard”), ex. 62 (varianta cu margine “soft”, folosind  $\nu$ -SVM);
- o condiție suficientă pentru ca variantele cu margine “hard” pentru cele două tipuri de probleme de optimizare *one-class SVM*, și anume *Max Margin* și *minimum enclosing ball* (MEB), în forma kernel-izată, să fie echivalente: ex. 31;
- deducerea *formeii duale* pentru problema  $\nu$ -SVM: ex. 32;
- deducerea *formeii duale* pentru problema SVR (*Support Vector Regression*), folosind funcție de cost / pierdere  $\varepsilon$ -senzitivă: ex. 33 (cu margine “hard”), ex. 64 (cu margine “soft” și (echivalent) cu funcție de cost  $\varepsilon$ -senzitivă); exemplificare / aplicare: ex. 63;
- teorema de reprezentare: ex. 65.