



© *M. Romanică*

Exerciții de învățare automată

Liviu Ciortuz, Alina Munteanu, Elena Bădăraș

2020 (varianta 2021f)

ISBN: 978-973-0-33148-6

0. Fundamente matematice:
probabilități, variabile aleatoare,
estimarea parametrilor unor distribuții probabiliste,
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. 83;
- calcularea unor *probabilități* elementare: ex. 2.a, ex. 80;
- calcularea unor *probabilități condiționate*: ex. 2.b, ex. 3, ex. 81.a;
- *regula de multiplicare*: ex. 83.b;
- *formula probabilității totale*: ex. 83.e;
formula probabilității totale – varianta condițională: ex. 86.cd;
- *independența* evenimentelor aleatoare: ex. 4, ex. 5, ex. 81.bc;
- *independența condițională* a evenimentelor aleatoare – legătura dintre forma „tare” și forma „slabă” a definiției pentru această noțiune: ex. 82;
- *formula lui Bayes* – aplicare: ex. 6, ex. 7, ex. 84, ex. 85;
formula lui Bayes – varianta condițională: ex. 86.b;
- recapitulare: probabilități elementare și probabilități condiționate – câteva proprietăți (A/F): ex. 8, ex. 86.

Variabile aleatoare

- funcție de *distribuție cumulativă* (engl., cumulative distribution function, c.d.f.), un exemplu: ex. 30;
- 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. 90;
- o *condiție suficientă* pentru ca $Var[X + Y] = Var[X] + Var[Y]$: independența variabilelor X și Y : ex. 21.c
- pentru *variabile aleatoare discrete*:
 - calcularea mediilor și a varianțelor – exemplificare: ex. 12, ex. 87, ex. 89.a;
 - definirea unei *variabile-indicator* cu ajutorul unui eveniment aleator; calcularea mediei acestei variabile: ex. 88;
 - regula de *înlanțuire* (pt. var. aleat.) – aplicare: ex. 13;
 - regula de *multiplicare* (pt. var. aleat.), varianta condițională – demonstrație: ex. 14;

- *independență*: ex. 89.b, ex. 91.a, ex. 16.a;
independență condițională: ex. 15, ex. 91.b, ex. 16.b, ex. 92.b;
o condiție suficientă pentru independența condițională: ex. 93;
- 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. 17.a, ex. 18, ex. 94;
 - dat fiind un p.d.f., să se calculeze o anumită probabilitate: ex. 17.c, ex. 95;
 - variabile aleatoare discrete vs. variabile aleatoare continue; p.m.f. vs. p.d.f.: ex. 97;
 - variabile aleatoare discrete și variabile aleatoare continue; independența și calculul de medii: ex. 98;
- *coeficientul de corelație* pentru două variabile aleatoare: ex. 19;
- *vector de variabile aleatoare*:
 - o proprietate: matricea de covarianță este simetrică și pozitiv definită: ex. 20; calculul matricei de covarianță când asupra vectorului de variabile aleatoare operăm transformări liniare: ex. 99;
- câteva inegalități de bază în teoria probabilităților: margini superioare pentru probabilități de forma $P(Z \geq t)$ și $P(Z - E[Z] \geq t)$: ex. 101;
- recapitulare: ex. 21, ex. 100.

Distribuții probabiliste uzuale

- distribuții discrete
 - distribuția *Bernoulli*:
 - suma de variabile identic distribuite; media sumei: ex. 22;
 - mixtură* de distribuții Bernoulli: ex. 104;
 - distribuția *binomială*:
 - verificarea condițiilor de definiție pentru p.m.f.: ex. 23.a;
 - calculul mediei și al varianței: ex. 23.b;
 - calcularea unor probabilități (simple și respectiv condiționale): ex. 102;
 - distribuția *categorială*:
 - calcularea unor probabilități și a unor medii: ex. 103;
 - mixtură* de distribuții categoriale: calculul mediei și al varianței: ex. 26;
 - distribuția *geometrică*:
 - calcularea numărului „așteptat” / mediu de „observații” necesare pentru ca un anumit eveniment să se producă: ex. 25;
 - distribuția *Poisson*:
 - verificarea condițiilor de definiție pentru p.m.f., calculul mediei și al varianței: ex. 24;

- **distribuții continue**
 - distribuția *continuă uniformă*:
exemplu de distribuție continuă uniformă unidimensională; calcularea mediei și a varianței: ex. 105;
calculul unei p.d.f. comune, pornind de la două variabile [urmând distribuții continue uniforme unidimensionale] independente; calculul unei anumite probabilități, folosind p.d.f. comună: ex. 27;
exemplu de distribuție continuă uniformă bidimensională; calcularea unei p.d.f. condiționale; verificarea independenței celor două distribuții marginale; calculul mediilor unor p.d.f.-uri condiționale: ex. 96, ex. 98;
 - distribuția *exponențială*:
verificarea condițiilor de definiție pentru p.d.f., calculul mediei și al varianței: ex. 28.a;
 - distribuția *Gamma*:
verificarea condițiilor de definiție pentru p.d.f., calculul mediei și al varianței: ex. 28.b;
 - distribuția *gaussiană unidimensională*:
verificarea condițiilor de definiție pentru p.d.f., calculul mediei și al varianței: ex. 29;
„standardizare“ (i.e., reducerea cazului nestandard la cazul standard): ex. 30;
 - distribuția *gaussiană bidimensională*: exemplificare; calcularea explicită a p.d.f.-ului, dat fiind vectorul de medii și matricea de covarianță: ex. 32;
o *proprietate* pentru distribuția *gaussiană bidimensională*: distribuția condițională a unei componente în raport cu cealaltă componentă este tot de tip gaussian; calculul parametrilor acestei distribuții condiționale: ex. 35;
 - distribuția *gaussiană multidimensională*:
matrice simetrice și pozitiv definite⁸⁹⁷ o proprietate de tip *factorizare* folosind *vectori proprii*: ex. 33;
densitatea distribuției gaussiene multidimensionale este într-adevăr p.d.f.:⁸⁹⁸ ex. 34;
o *proprietate* importantă, în cazul în care matricea de covarianță este diagonală: p.d.f.-ul comun este produsul p.d.f.-urilor marginale (care sunt independente): ex. 31;
 - În ce privește specificul datelor generate de distribuții gaussiene multidimensionale:⁸⁹⁹
 - în cazul cel mai general (deci când matricea Σ nu este neapărat diagonală), datele generate de acest tip de distribuție tind să se grupeze în elipse (corpuri elipsoidale) cu *axele de simetrie* [desigur, perpendiculare, dar] în general ne-paralele cu *axele de coordonate*;
 - dacă matricea de covarianță Σ este diagonală, atunci datele generate tind să se grupeze în elipse (dacă se lucrează în \mathbb{R}^2) sau, mai general, corpuri elipsoidale având axele de simetrie paralele cu axele sistemului de coordonate;
 - 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;

⁸⁹⁷ Așa sunt matricele de covarianță ale variabilelor gaussiene multidimensionale.

⁸⁹⁸ Adică satisface condițiile din definiția noțiunii de p.d.f.

⁸⁹⁹ Vedeti corespondența imediată cu alura curbilor de izocontur / nivel determinate de aceste distribuții.

- *mixturi de distribuții gaussiane multidimensionale*: exprimarea vectorului de medii și a matricei de covarianță în funcție de mediile și matricele de covarianță ale distribuțiilor componente: ex. 108;
- *mixturi de distribuții oarecare*: calculul mediilor și al varianțelor (în funcție de mediile și matricele de covarianță ale distribuțiilor componente): ex. 109;
în cazul unei distribuții multidimensionale de tip mixtură $p(x) = \sum_{k=1}^K \pi_k p(x|z_k)$, dacă partiționează variabila $x \in \mathbb{R}^d$ în două, sub forma $x \stackrel{not.}{=} (x_1, x_2)$, atunci variabila condiționată $x_2|x_1$ urmează tot o distribuție de tip mixtură: ex. 36;
- distribuția Bernoulli și distribuția normală standard: *intervale de încredere, legea numerelor mari, teorema limită centrală*; aplicație la calculul erorii reale a unui clasificator: ex. 37;
- funcția generatoare de momente pentru o variabilă aleatoare: ex. 110;
- *familia de distribuții exponențiale*: ex. 38, ex. 112, precum și ex. 40.a și ex. 41.a de la capitolul *Metode de regresie*;
- chestiuni recapitulative (corespondența dintre nume de distribuții și expresiile unor p.d.f.-uri date): ex. 111.

Estimarea parametrilor unor distribuții probabiliste uzuale

- distribuția Bernoulli: ex. 40 (+MAP, folosind distr. Beta), ex. 39, ex. 115 (un caz particular), ex. 116.a, ex. 114 (bias-ul și varianța estimatorului MLE);
- distribuția binomială: ex. 116.bc;
- distribuția categorială: ex. 117, ex. 118 (+MAP, folosind distr. Dirichlet);
- distribuția geometrică: ex. 119 (+MAP, folosind distr. Beta);
- distribuția Poisson: ex. 41 (+MAP, folosind distr. Gamma);
- distribuția uniformă continuă: calcul de probabilități și MLE: în \mathbb{R} : ex. 42, ex. 120, ex. 121; în \mathbb{R}^2 : ex. 43, ex. 122;
- distribuția gaussiană unidimensională: MLE pt. μ , considerând σ^2 cunoscut: ex. 45 (+MAP, folosind distribuția gaussiană); MLE pt. σ^2 , atunci când nu se impun restricții asupra lui μ : ex. 46 (+deplasare); MLE pt. σ^2 , atunci când $\mu = 0$: ex. 124 (+nedeplasare);
- distribuția gaussiană multidimensională: ex. 48, ex. 125 (+MAP, folosind distr. Gauss-Wishart);
- distribuția exponențială: ex. 44, ex. 123 (+MAP, folosind distr. Gamma);
- distribuția Gamma: ex. 47 și ex. 126 (ultimul, folosind metoda gradientului și metoda lui Newton);
- existența și unicitatea MLE: ex. 49;
- MLE și parametrizarea alternativă: ex. 127.

Elemente de teoria informației

- *definiții și proprietăți imediate* pentru entropie, entropie comună, entropie condițională specifică, entropie condițională medie, câștig de informație: ex. 50; *exemplificarea* acestor noțiuni (varianta discretă): ex. 51, ex. 52, ex. 130; *exemple* de calculare a entropiei unor variabile aleatoare continue: distribuția continuă uniformă (ex. 129.a), distribuția exponențială (ex. 55.a), distribuția Gamma (ex. 55.b), distribuția gaussiană unidimensională (ex. 129.b) și distribuția gaussiană multidimensională (ex. 129.c);
o *proprietate* a entropiei: nenegativitatea: ex. 50.a, ex. 138.a;
o *margine superioară* pentru valoarea entropiei unei variabile aleatoare discrete: ex. 131;
două *proprietăți* ale câștigului de informație: nenegativitatea (ex. 57.c, ex. 134) și $IG(X; Y) = 0 \Leftrightarrow$ variabila X este independentă de Y (ex. 57.c);
- re-descoperirea definiției entropiei, pornind de la un set de *proprietăți dezirabile*: ex. 56;
- o *aplicație* pentru câștigul de informație: *selecția de trăsături*: ex. 53;
- *entropia comună*: exemplu de calculare: ex. 130.d
forma particulară a relației de „înlănțuire” în cazul variabilelor aleatoare independente: ex. 54, ex. 132, ex. 138.c;
entropie comună și condiționată: formula de „înlănțuire” condițională: ex. 133;
- *entropia relativă* (divergența Kullback-Leibler / KL):
definiție și proprietăți elementare: ex. 57.a;
exprimarea câștigului de informație cu ajutorul entropiei relative: ex. 57.b;
o proprietate de tip „relație de înlănțuire”: ex. 136;
echivalența dintre minimizarea entropiei relative și maximizarea funcției de log-verosimilitate: ex. 137;
- *cross-entropie*: definiție, o proprietate (nenegativitatea) și un exemplu simplu de calculare a valorii cross-entropiei: ex. 58;
un exemplu de *aplicație* pentru cross-entropie: selecția modelelor probabiliste: ex. 135;
- *inegalitatea lui Gibbs*: un caz particular; comparație între valorile entropiei și ale cross-entropiei: ex. 59;
- entropia văzută ca o *funcțională* în raport cu p ; calcululul *derivatei funcționale* a entropiei în raport cu p : ex. 60.

Funcții-nucleu

- aflarea funcției de „mapare” a trăsăturilor care corespunde unei funcții-nucleu date: ex. 61, ex. 139, ex. 140.a;
comparații asupra numărului de operații efectuate la calcularea valorii unor funcții-nucleu (în spațiul inițial vs. spațiul nou de „trăsături”): ex. 140.b;
calculul distanțelor euclidiene în spații de „trăsături” folosind doar funcții-nucleu: ex. 65;
- *teorema lui Mercer* (1909): condiții necesare și suficiente pentru ca o funcție să fie funcție-nucleu: ex. 62.ab;

- rezultate de tip „constructiv“ pentru [obținerea de noi] funcții-nucleu: ex. 62.c, 63, ex. 64, ex. 142, ex. 70.b; contraexemple: ex. 70.ac, ex. 143; „normalizarea“ funcțiilor-nucleu: ex. 141;
- 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. 144;
- un exemplu de funcție-nucleu care servește la a măsura similaritatea dintre două imagini oarecare: ex. 66;
- o funcție-nucleu particulară, care asigură separabilitate liniară [în spațiul de trăsături] pentru orice set de instanțe de antrenament: ex. 145;
- funcția-nucleu gaussiană / *funcția cu baza radială* (engl., Radial Basis Function, RBF):
 - demonstrația faptului că RBF este într-adevăr funcție-nucleu: ex. 67;
 - funcția de „mapare“ corespunzătoare funcției-nucleu RBF la valori într-un spațiu [de „trăsături“] de dimensiune infinită: ex. 68;
 - proprietăți simple ale nucleului RBF: ex. 69, ex. 146;
 - 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. 148.

Metode de optimizare în învățarea automată

- (P0) definiții, caracterizări și câteva proprietăți pentru funcții convexe: ex. 71;
- inegalitatea lui Jensen: ex. 131;
- *metoda gradientului*, exemplificare: ex. 73.c, ex. 151, ex. 152; implementare: ex. 149.
metoda lui Newton, exemplificare: ex. 73.d;
- (P1) condiții suficiente pentru convergența metodei gradientului: ex. 150;
- (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. 153;
- (P3) *reparametrizarea liniară* a atributelor nu afectează [rezultatele obținute cu] metoda lui Newton, însă afectează metoda gradientului: ex. 154;
- *metoda dualității Lagrange*:
 - (P4) demonstrarea proprietății de *dualitate slabă*: ex. 74;
 - (P5) demonstrarea unei părți din *teorema KKT*: ex. 75;
 - exemple de aplicare: ex. 76, ex. 77, ex. 78, ex. 79, ex. 131, ex. 155;
 - un exemplu de problemă de optimizare convexă pentru care condițiile KKT nu sunt satisfăcute: ex. 156;
 - aplicarea metodei dualității Lagrange la determinarea distribuțiilor unidimensionale care — în anumite condiții — au entropii maxime: ex. 157;

- două variante a algoritmului Perceptron,⁹⁰⁰ pentru care relația de actualizare a ponderilor se obține rezolvând [câte] o problemă de optimizare [convexă] cu restricții.
- metoda descreșterii pe coordonate, în contextul regresiei liniare cu regularizare L_1 (regresia Lasso): ex. 11 de la capitolul *Metode de regresie*; metoda subgradientului, tot în contextul regresiei liniare cu regularizare L_1 (regresia Lasso): ex. 28 de la capitolul *Metode de regresie*.

⁹⁰⁰Vedeți problema 16 de la capitolul *Rețele neuronale artificiale*.

1. Metode de regresie

Sumar

Noțiuni preliminare

- estimarea parametrilor unor distribuții probabiliste uzuale (în special distribuția Bernoulli, distribuția gaussiană, distribuția Laplace); vedeți secțiunea corespunzătoare de la cap. *Fundamente*;
- elemente de calcul vectorial (în particular, produsul scalar) și de calcul matriceal: ex. 33 de la cap. *Fundamente*; norma L_2 (euclidiană) și norma L_1 : ex. 3, ex. 26, ex. 11; calculul derivatelor parțiale [de ordinul întâi și al doilea]: ex. 7; reguli de derivare cu argumente vectoriale: ex. 1;
- metode de optimizare (în speță pentru aflarea maximului / minimului unei funcții reale, derivabile): metoda analitică, metoda gradientului, metoda lui Newton; exemplificare: ex. 73 de la capitolul de Fundamente.

Regresia liniară

- prezentarea *generală* a metodei regresiei liniare:⁹⁰¹
 - MLE și corespondența cu estimarea în sens LSE (least squared errors): ex. 3.A.abc; rezolvare folosind matricea de design (X'): ex. 3.A.d; particularizare pentru cazul unidimensional: ex. 1.ab, ex. 20; exemplificare pentru cazul unidimensional (ex. 2, ex. 21) și pentru cazul bidimensional (ex. 23.a, ex. 32);
 - (P1) *scalarea atributelor* nu schimbă predicțiile obținute (pentru instanțele de test) cu ajutorul *formulelor analitice*: ex. 4, 39.a;
 - (P2) adăugarea de noi trăsături / attribute nu mărește suma pătratelor erorilor: ex. 25;
 - 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. 7;
 - MAP și corespondența cu *regularizarea* de normă L_2 (regresia *ridge*): ex. 3.C; particularizare pentru cazul unidimensional: ex. 1.c;
 - (P3) efectul de diminuare a ponderilor (engl., weight decay) în cazul regularizării de normă L_2 (respectiv L_1 — regresia *Lasso*) a regresiei liniare, în comparație cu cazul neregularizat: ex 26.b (respectiv ex. 26.a);
- bias-ul și [co]varianța estimatorului regresiei liniare; bias-ul regresiei *ridge*: ex. 5;
- regresia polinomială [LC: mai general: folosirea așa-numitelor *funcții de bază*]: ex. 3.B; exemplificare pentru cazul bidimensional: CMU, 2015 spring, T. Mitchell, N. Balcan, HW4, pr. 1;

⁹⁰¹În mod implicit, în această secțiune se va considera că termenul-zgomot este modelat cu distribuția gaussiană (dacă nu se specifică altfel, în mod explicit).

- cazul regresiei liniare cu termen de regularizare L_2 (regresia *ridge*): deducerea regulilor de actualizare pentru *metoda gradientului ascendent*: varianta “batch” / “steepest descent”: ex. 6.a; și varianta stohastică / secvențială / “online”: ex. 6.b; exemplu de aplicare: ex. 24;
- cazul regresiei liniare cu termen de regularizare L_1 (regresia Lasso): rezolvare cu *metoda descenderii pe coordonate* (engl., “coordinate descent”): ex. 11; rezolvare cu *metoda subgradientului* (aplicare la selecția de trăsături): ex. 28;
- regresia liniară în cazul zgomotului modelat cu distribuția *Laplace* (în locul zgomotului gaussian): ex. 8.B; exemplificare pentru cazul bidimensional: ex. 23.c; rezolvare în cazul unidimensional [chiar particularizat] cu ajutorul derivatei, acolo unde aceasta există: ex. 29;
- regresia liniară și *overfitting-ul*: ex. 12;
- regresie liniară folosită pentru *clasificare*: exemplificare: ex. 32;
- cazul *multivaluat* al regresiei liniare, reducerea la cazul uninomial: ex. 31;
- regresia liniară cu regularizare L_2 (regresia *ridge*), *kernel-izarea* ecuațiilor „normale”: ex. 9; aplicare: ex. 27; (P4) folosind nucleu RBF, eroarea la antrenare devine 0 atunci când parametrul de regularizare λ tinde la 0: ex. 10;
- *regresia liniară ponderată*: ex. 8.A; particularizare pentru cazul unidimensional [neparametric]: CMU, 2010 fall, Aarti Singh, midterm, pr. 4; particularizare / exemplificare pentru cazul bidimensional: ex. 23.b; cazul multivaluat, cu regularizare L_2 : Stanford, 2015 fall, Andrew Ng, midterm, pr. 2; o *proprietate* a regresiei liniare local-ponderate [demonstrată în cazul unidimensional]: „netezirea” liniară: ex. 30.

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. 13, 39.b;
 - particularizare* pentru cazul datelor din \mathbb{R}^2 : ex. 33 (inclusiv *regularizare* L_1 / estimarea parametrilor în sens MAP, folosind o distribuție a priori Laplace);
- (P0) *granița de decizie* pentru regresia logistică: ex. 33.d;
- (P1) funcția de log-verosimilitate în cazul regresiei logistice este concavă (decare un maxim global), fiindcă matricea hessiană este pozitiv definită: ex. 14;

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. 34;

- (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 multidimensională sferică —, în comparație cu cazul estimării parametrilor în sensul MLE: ex. 15;
- 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. 16;
 - regresia logistică kernel-izată: adaptarea metodei gradientului: ex. 17;
 - regresia logistică n -ară (așa-numita regresie *softmax*): calculul funcției de log-verosimilitate, cu regularizare L_2 , deducerea regulilor de actualizare a ponderilor, folosind metoda gradientului: ex. 18;
- (P4) echivalența cu un anumit tip de mixtură de distribuții gaussiene multidimensionale: ex. 36;
- (P5) o [interesantă] proprietate comună pentru regresia liniară și regresia logistică: ex. 35;
- 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. 39.c;
- comparații între regresia logistică și alți clasificatori (Bayes Naiv, ID3): ex. 33.c, ex. 37.ab;
- (P6) teorema de reprezentare: ex. 19, ex. 38.

Modele liniare generalizate (GLM)

- condiții suficiente pentru *concavitatea* funcției de log-verosimilitate: ex. 42;
- particularizare pentru cazul distribuției geometrice: ex. 40;
- particularizare pentru cazul distribuției gaussiene unidimensionale: ex. 41.

2. Clasificare bayesiană

Sumar

Noțiuni preliminare

- probabilități și probabilități condiționate;
- formula lui Bayes: ex. 5.b;
cap. *Fundamente*, ex. 6, ex. 7, ex. 84, ex. 85;
- independența [condițională a] evenimentelor aleatoare:
cap. *Fundamente*, ex. 4, ex. 81, ex. 82;
- independența [condițională a] variabilelor aleatoare: ex. 9, ex. 10, ex. 12,
ex. 31-38; vedeți și cap. *Fundamente*, ex. 15, ex. 27, ex. 89.b, ex. 98, ex. 96;
- distribuții probabiliste comune, marginale și condiționale: ex. 8, ex. 10, ex. 12,
ex. 31; vedeți și cap. *Fundamente*, ex. 13, ex. 14;
- distribuția gaussiană: de la cap. *Fundamente*, ex. 29, ex. 30 (pentru cazul
unidimensional), ex. 32 (pentru cazul bidimensional), ex. 20, ex. 31, ex. 33,
ex. 34 (pentru cazul multidimensional);
- estimarea parametrilor pentru distribuții de tip Bernoulli, categorial și gaus-
sian (ultimul doar pentru cazul clasificării bayesiene de tip gaussian);⁹⁰²
- ipoteze MAP vs. ipoteze ML:
formulare [ca soluții la] probleme de optimizare:⁹⁰³ ex. 25;
exemplificare: ex. 1, ex. 2, ex. 3, ex. 24, ex. 37;
exemplificare în cazul arborilor de decizie: ex. 4;
- regresia logistică, chestiuni introductive:⁹⁰⁴ de la cap. *Metode de regresie*,
ex. 13.

Algoritmi de clasificare bayesiană

- Algoritmul Bayes Naiv și algoritmul Bayes Optimal:⁹⁰⁵
formulare ca probleme de optimizare / estimare în sens MAP: cartea ML,
pag. 157;
pseudo-codul algoritmului Bayes Naiv (cf. cartea ML, pag. 177):

Training:

```

for each value  $v_j$  of the output attribute
   $\hat{P}(v_j) \leftarrow \text{estimate } P(v_j)$ 
  for each value  $a_i$  of each input attribute  $a$ 
     $\hat{P}(a_i|v_j) \leftarrow \text{estimate } P(a_i|v_j)$ 

```

⁹⁰²De la cap. *Fundamente*, pentru estimarea parametrului unei distribuții Bernoulli vedeți ex. 40 și ex. 114.a, pentru estimarea parametrilor unei distribuții categoriale vedeți ex. 117, iar pentru estimarea parametrilor unei distribuții gaussiene vedeți ex. 45, ex. 46, ex. 124 (pentru cazul unidimensional) și ex. 48 (pentru cazul multidimensional).

⁹⁰³Vedeți cartea ML, pag. 156-157.

⁹⁰⁴Vedeți draftul capitolului suplimentar pentru cartea ML a lui T. Mitchell, *Generative and discriminative classifiers: Naive Bayes and logistic regression* (în special secțiunea 3).

⁹⁰⁵La 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.

Classification of the test instance (a_1, a_2, \dots, a_n) :

$$v_{NB} = \operatorname{argmax}_{v_j \in V} \hat{P}(v_j) \prod_{a_i \in x} \hat{P}(a_i | v_j).$$

Justificarea regulii de decizie pentru algoritmul Bayes Naiv:⁹⁰⁶

$$\begin{aligned} v_{MAP} &= \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2 \dots a_n) = \operatorname{argmax}_{v_j \in V} \frac{P(a_1, a_2 \dots a_n | v_j) P(v_j)}{P(a_1, a_2 \dots a_n)} \\ &= \operatorname{argmax}_{v_j \in V} P(a_1, a_2 \dots a_n | v_j) P(v_j) \stackrel{\text{indep. cdt.}}{=} \operatorname{argmax}_{v_j \in V} \prod_i P(a_i | v_j) P(v_j) \stackrel{\text{not.}}{=} v_{NB}. \end{aligned}$$

- exemple de aplicare: ex. 5, ex. 7, ex. 8, ex. 9, ex. 26, ex. 27, ex. 28;
- aplicarea / adaptarea algoritmului Bayes Naiv pentru clasificare de texte:⁹⁰⁷ ex. 6, ex. 29;
folosirea regulii “add-one” [a lui Laplace] pentru „netezirea” parametrilor: ex. 6, ex. 30;
- calculul ratei medii a erorilor pentru algoritmi Bayes Naiv și Bayes Optimal: ex. 10, ex. 11, ex. 31, ex. 32, ex. 33, ex. 34, ex. 38;
- evidențierea grafică a neconcordanței predicțiilor făcute de clasificatorii Bayes Naiv și Bayes Optimal: ex. 12.

Proprietăți ale algoritmilor Bayes Naiv și Bayes Optimal

- (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 Optimal) în faza de testare coincid;
- (P1) numărul de parametri necesari de estimat din date: liniar pentru Bayes Naiv $(2d + 1)$ și exponențial pentru Bayes Optimal $(2^{d+1} - 1)$:⁹⁰⁸ ex. 7.e, ex. 28.ab, ex. 33.ac;
- (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 Optimal 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);

⁹⁰⁶ Justificarea regulii de decizie pentru algoritmul Bayes Optimal / Comun:

$$\begin{aligned} v_{MAP} &= \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2 \dots a_n) = \dots \\ &= \operatorname{argmax}_{v_j \in V} P(a_1, a_2 \dots a_n | v_j) P(v_j) = \operatorname{argmax}_{v_j \in V} P(a_1, a_2 \dots a_n, v_j) \stackrel{\text{not.}}{=} v_{JB}. \end{aligned}$$

⁹⁰⁷ 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.

⁹⁰⁸ Numărul de parametri indicați în paranteze se referă la cazul când atât atributele de intrare cât și atributul de ieșire sunt de tip Bernoulli.

- (P4) acuratețea [la clasificare a] algoritmului Bayes Naiv scade atunci când unul sau mai multe atribute de intrare sunt duplicate: ex. 10.d, ex. 31.def;
- (P5) în cazul „învățării” unei funcții booleene (oarecare), rata medie a erorii produse la antrenare de către algoritmul Bayes Optimal (spre deosebire de Bayes Naiv!) este 0: ex. 33.d;
- (P6) complexitatea de eșantionare: de ordin logaritmă pentru Bayes Naiv și de ordin exponențial pentru Bayes Optimal: 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. 36, ex. 38.

Algoritmi Bayes Naiv și Bayes Optimal cu variabile de intrare *de tip gaussian*

- Aplicație: G[N]B: ex. 15, ex. 40 și ex. 48; GJB: ex. 45, ex. 46 și ex. 47; GNB vs. GJB: ex. 21.
- Numărul de parametri necesari de estimat din date: ex. 42.
- Proprietăți:
 - (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: demonstrație: ex. 17; exemplificare în \mathbb{R} : ex. 40.a; exemplificare în \mathbb{R}^2 : ex. 41.c;
 - (P1') similar, presupunem că variabila de ieșire este booleană; dacă variabilele de intrare (notație: $X = (X_1, \dots, X_d)$) au distribuțiile [comune] condiționale $X|Y = 0$ și $X|Y = 1$ de tip gaussian [multidimensional], 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, ex. 20.a.i – ii;
 - (P2') când variabilele de intrare satisfac condiții mixte de tip (P0') sau (P7), atunci concluzia – separare liniară – se menține: ex. 44.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. 22.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. 22.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. 23.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 (în mod necesar) forma regresiei logistice: ex. 40.bc, ex. 43;
similar, pentru algoritmul Bayes Optimal gaussian, atunci când $\Sigma_0 \neq \Sigma_1$, ecuația separatorului decizional este de ordin pătratic: ex. 19, ex. 20.a.iii – vi, ex. 42.e (separatorul decizional este un cerc), ex. 42.f (o hiperbolă), ex 47 (o reuniune de două drepte);
- (P6') parametrii algoritmilor Bayes Naiv gaussian și Bayes Optimal gaussian se pot estima în timp liniar în raport cu numărul de instanțe din setul de date de antrenament: ex. 41.bd, ex. 23.b.

3. 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 (cf. cartea ML, pag. 232):

Training:

Store all training examples.

Classification:

Given a query/test instance x_q ,
first locate the k nearest training examples x_1, \dots, x_k ,
then take a vote among its k nearest neighbors

$$\operatorname{argmax}_{v \in V} \sum_{i=1}^k \delta(v, f(x_i))$$

where $\delta(a, b) = 1$ if $a = b$, and $\delta(a, b) = 0$ if $a \neq b$.

- 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;
- 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 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 neparametrică; alegerea lui k : CV: ex. 23.c;

- CVLOO: ex. 3, ex. 12.b, ex. 15.bc, 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 atribute 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 Optimal: ex. 10, ex. 25.

Comparații cu alți algoritmi de clasificare automată

- ID3: ex.11.b, ex.12.c, ex. 13, ex. 24.b;
- SVM: ex.12.d, 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);⁹⁰⁹
- 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.

⁹⁰⁹Sectiunea 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.

4. Arbori de decizie

Sumar

Noțiuni preliminare

- partiție a unei mulțimi: ex. 83 de la cap. *Fundamente*;
- proprietăți elementare ale funcției logaritm; formule uzuale pentru calcule cu logaritmi;
- Elemente de *teoria informației* (vedeți secțiunea corespunzătoare din cap. *Fundamente*):
 - entropie, definiție: T. Mitchell, *Machine Learning*, 1997 (desemnată în continuare simplu prin *cartea ML*), pag. 57; ex. 2.a, ex. 39.a, ex. 35.a;
 - entropie condițională specifică: ex. 14.a;
 - entropie condițională medie: ex. 2.cd, ex. 35.c;
 - câștig de informație (definiție: *cartea ML*, pag. 58): ex. 2.cd, ex. 5.a, ex. 33, ex. 39.b, ex. 35.e;
- *arbori de decizie*, văzuți ca *structură de date*: ex. 1, ex. 31
și, respectiv, ca program în logica propozițiilor: ex. 2.e, ex. 38.bc;
- (P0) *expresivitatea arborilor de decizie* cu privire la *funcții boolene*: ex. 32;
- *spațiu de versiuni* pentru un concept (de învățat): ex. 1, ex. 31, ex. 37.

Algoritmul ID3

- pseudo-cod, versiune simplificată:


```

create the root node;
assign all training examples to the root;

Main loop:
1.  $A \leftarrow$  the “best” decision attribute for the current node;
2. for each value of the attribute  $A$ , create a new descendant of the node;
3. sort the training examples to the leaf nodes;
4. if the training examples are perfectly classified, then STOP;
   else iterate over the new leaf nodes;
      
```

pseudo-cod, versiune completă: *cartea ML*, pag. 56;
- *bias-ul inductiv*: *ibidem*, pag. 63-64;
- exemple simple de aplicare: ex. 2, ex. 3, ex. 4.a, ex. 5, ex. 36, ex. 37.a, ex. 38, ex. 39, ex. 40;
- 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. 31, ex. 37;
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 / testare și *iii.* în final să aibă o [cât mai] bună putere de *generalizare*;⁹¹⁰
- ID3 ar putea fi văzut și ca algoritm de *optimizare*;⁹¹¹
- *greedy* \Rightarrow nu garantează obținerea soluției optime d.p.v. al numărului de niveluri / noduri:
ex. 4, ex. 21.a, ex. 38 (vs. ex. 37.b, ex. 3.b), ex. 46;
- de tip *divide-et-impera* (\Rightarrow “*Iterative Dichotomizer*”), recursiv;
- *1-step look-ahead*;
- complexitate de timp, cf. *Weka book*;⁹¹²
la antrenare, în anumite condiții: $\mathcal{O}(dm \log m)$; la testare $\mathcal{O}(d)$,
unde d este numărul de attribute, 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] *arborele* 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. 42;⁹¹³ [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. 45.bc;⁹¹⁴
 - *robustețe la „zgomote” și overfitting*: ex. 10, ex. 21.bc, ex. 45, ex. 67.b;⁹¹⁵
 - *zone de decizie și granițe de separare / decizie* pentru arbori de decizie cu variabile continue: ex. 10, ex. 44, ex. 45, ex. 46.
Observație: Zonele de decizie produse de algoritmul ID3 nu sunt în mod neapărat unice, fiindcă arborele de decizie creat de ID3 nu este determinat în mod unic (vedeți proprietatea (P2), mai jos).

Extensii / variante ale algoritmului ID3

- *attribute cu valori continue*: ex. 10-12, ex. 14.c, ex. 43-47; cap. *Învățare bazată pe memorare*, ex. 11.b;
- *attribute discrete cu multe valori*: ex. 13;
- *attribute cu valori nespecificate / absente pentru unele instanțe*;
- *attribute cu diferite costuri asociate*: ex. 14;

⁹¹⁰LC: 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 attributele de *intrare*, care sunt adăugate la model câte unul pe rând.

⁹¹¹LC: 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*, și pentru care entropia să fie *minimală* (vedeți proprietatea de structuralitate de la ex. 56 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.

⁹¹²“Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations”, Ian Witten, Eibe Frank (3rd ed.), Morgan Kaufmann Publishers, 2011.

⁹¹³De asemenea, ex. 4.ab de la capitolul *Clasificare bayesiană*.

⁹¹⁴De asemenea, ex. 4.ab de la capitolul *Clasificare bayesiană*.

⁹¹⁵De asemenea, ex. 4.ab de la capitolul *Clasificare bayesiană*.

- reducerea caracterului “eager” al învățării: ex. 16;
- reducerea caracterului “greedy” al învățării:
IG cu “2-step look-ahead”: ex. 17, ex. 18;
variante de tip “look-ahead” specifice atributelor continue: ex. 48;
- folosirea altor *măsuri de „impuritate”* în locul câștigului de informație:
Gini Impurity, Misclassification Impurity: ex. 15;
- reducerea *overfitting*-ului:
reduced-error pruning (folosind un set de date de validare): cartea ML, pag. 69-71; A. Cornuéjols, L. Miclet, *Apprentissage artificiel*, 2010, pag. 418-421;
rule post-pruning: cartea ML, pag. 71-72; ex. 50
top-down vs. bottom-up pruning, folosind un criteriu bazat pe câștigul de informație: ex. 19, ex. 49;
pruning folosind testul statistic χ^2 : ex. 20, 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. 37-38;
- (P2) arborele produs de algoritmul ID3 *nu* este în mod neapărat *unic*: ex. 3, ex. 37;
- (P3) arborele ID3 nu este neapărat *optimal* (ca nr. de noduri / niveluri): ex. 4, ex. 21, ex. 38;
- (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. 41;
- (P7) *granițele de separare / decizie* pentru arborii ID3 cu *attribute de intrare continue* sunt întotdeauna paralele cu axele de coordonate: ex. 10, ex. 12, ex. 44, ex. 45, ex. 46 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 attributele de intrare sunt categoriale: numărul de attribute: ex. 52.c;
- (P9) o *margină superioară* pentru *adâncimea* unui arbore de decizie când attributele 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 attribute de intrare, atunci când acestea (attributele de intrare) sunt binare: ex. 9.
- Alte metode de învățare automată bazate pe arbori: arbori de regresie (CART).

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. 63), funcții de cost / pierdere (engl., loss function) (vedeți ex. 29 și ex. 23.b);
- pseudo-codul algoritmului AdaBoost: ex. 22;⁹¹⁶
 - proprietăți de bază: ex. 22;
 - alte proprietăți, precum și două margini superioare pentru eroarea la antrenare: ex. 23.a-d;
 - convergența algoritmului, precum și o condiție suficientă pentru învățabilitate empirică γ -slabă: ex. 23.e;
- exemple de aplicare: ex. 24, 54, 55, 56, 57, 58.
- 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 greedy (dacă la fiecare iterație se alege cea mai bună ipoteză „slabă“), algoritm de optimizare secvențială (minimizează o margine superioară pentru eroarea la antrenare);
- *învățabilitate empirică γ -slabă*:
 - definiție: ex. 23.e
 - exemplificarea unor cazuri când nu există *garanție* pentru învățabilitate γ -slabă: ex. 25, 60;
- AdaBoost ca algoritm de *optimizare secvențială* în raport cu funcția de cost / „pierdere“ negativ-exponențială: ex. 26;
- marginea de votare: ex. 27, 28 și 61;
 - proprietăți ale marginilor de votare:
 - $\text{Margin}_k(i) \in [-1, +1]$
 - $\text{Margin}_k(x_i) = y_i \bar{f}_k(x_i)$ unde $\bar{f}_k(x_i) \stackrel{\text{not.}}{=} \sum_{t=1}^k \bar{\alpha}_t h_t(x_i)$
 - x_i este corect clasificat la iterația $k \Leftrightarrow \text{Margin}_k(i) \geq 0$
 - $\text{Margin}_k(x_i) > \text{Margin}_k(x_j) \Leftrightarrow D_{k+1}(i) < D_{k+1}(j)$;
- *selectarea trăsăturilor* folosind AdaBoost; aplicare la clasificarea de documente: ex. 62;
- o variantă generalizată a algoritmului AdaBoost:⁹¹⁷ ex. 29 și ex. 65;
- recapitulare (întrebări cu răspuns *adevărat* / *fals*): ex. 30 și 66.

⁹¹⁶Vedeți pseudo-codul algoritmului AdaBoost.M1 din cartea „Probabilistic machine learning: An introduction“ de Kevin Murphy, pag. 577, MIT Press, 20201. (Cf. <https://probml.github.io/pml-book/book1.html>, accesat la 19.02.2021.)

⁹¹⁷Pentru AdaBoost văzut ca instanță a unui algoritm general de învățare ansamblistă bazat pe minimizarea secvențială a unei funcții de cost / pierdere, vedeți ex. 26.d.

- **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. 22.vii);

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. 22.ii) și

$\varepsilon_t \in (0, 1/2) \Rightarrow Z_t \in (0, 1)$ (ex. 22.iii).

(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.a).

Produsul $y_i f_t(x_i)$ se numește *margină 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.b);

(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 α_t , ceea ce conduce la $\alpha_t = \ln \sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}}$

(ex. 23.c);

(P5') $\varepsilon_i > \varepsilon_j \Leftrightarrow \alpha_i < \alpha_j$ (consecință imediată din relația $\alpha_t = \ln \sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}}$): ex. 22.v;

(P6) O consecință din relația (193) și (P5): $D_{t+1}(i) = \begin{cases} \frac{1}{2\varepsilon_t} D_t(i), & i \in M \\ \frac{1}{2(1 - \varepsilon_t)} D_t(i), & i \in C \end{cases}$

(ex. 22.iv);

(P7) $err_S(H_t)$ nu neapărat descrește de la o iterație la alta; în schimb, descresc marginile sale superioare: $\prod_{t'=1}^t Z_{t'}$ și $\exp(-\sum_{t'=1}^t \gamma_{t'}^2)$ (ex. 23.be);

(P8) O condiție suficientă pentru învățabilitate γ -slabă, bazată pe marginea de votare: la fiecare iterație a algoritmului AdaBoost, media marginilor de votare ale instanțelor de antrenament în raport cu distribuția D_t să fie de cel puțin 2γ (ex. 28);

(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 m compași de decizie (ex. 63.a);

(P9') Orice mulțime de instanțe distincte [și etichetate] din \mathbb{R} este γ -slab învățabilă cu ajutorul compașilor de decizie (ex. 64).

- Alte metode de învățare ansamblistă bazate pe arbori de decizie: Bagging, Random Forests.

5. 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. 35;
- 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. 36, ex. 38;
- distanța (cu sau fără semn) de la un punct la o dreaptă (respectiv la un plan, sau mai general la un hiperplan): ex. 1, ex. 5, ex. 6.c;
- 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. 35, ex. 75, ex. 76, ex. 77, ex. 155, ex. 156 de la capitolul de *Fundamente*;
- *separabilitate* liniară, separator optimal, *margină* geometrică: ex. 38, 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;⁹¹⁸
- (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. 36, ex. 38, ex. 39 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. 37;
- exemple de [funcții de] *mapare a atributelor*, cu scopul de a obține separabilitate liniară în spațiul de trăsături: ex. 6.d, ex. 8, ex. 9, ex. 40.bd, ex. 41.a; rezolvarea directă a problemei SVM primale în [noul] spațiu de trăsături; identificarea separatorului neliniar din spațiul inițial [de trăsături]: ex. 9.de, ex. 40.ce, ex. 41.b-e;
- (•) deducerea *forme duale* pentru problema SVM cu margine “hard”: ex. 10;
- *exemplificarea* a două modalități de găsire a soluției forme 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 Φ dată: prin optimizare directă (ex. 11, ex. 42), respectiv prin folosirea relațiilor de legătură cu soluția problemei primale: ex. 12;

⁹¹⁸Vedeți și Andrew Ng (Stanford), Lecture Notes, part V, section 3.

- (P1) dacă datele de antrenament sunt liniar separabile, atunci problema de optimizare SVM în formă primală are soluție unică, întrucât funcția obiectiv a acestei probleme este strict convexă: ex. 49.a (vedeți *Observația*).
- (P2) efectul *multiplicării valorilor atributelor* cu o constantă pozitivă asupra separatorului obținut de SVM), atunci când datele de antrenament sunt liniar separabile cu ajutorul unui hiperplan care trece prin originea sistemului de coordonate: ex. 27.a;⁹¹⁹
- (P3) 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. 49, ex. 57;
- vedeți proprietățile (P6) și (P8) enunțate mai jos (la secțiunea despre C-SVM), care sunt valabile și în cazul SVM [cu margine “hard”];
- *comparații* între SVM [având, eventual, diferite funcții-nucleu] și alți clasificatori: ex. 45, ex. 46, ex. 57; vedeți și ex. 12 de la capitolul *Învățare bazată pe memorare*.

SVM cu margine “soft” (C-SVM):

- (•) deducerea *formeii duale* pentru problema SVM cu margine “soft” (C-SVM): ex. 13.a-d;
- (P4) legătura dintre valorile (și intervalele de valori) pentru multiplicatorii Lagrange α_i pe de o parte, și valorile (și intervalele de valori) pentru *marginile funcționale* $y_i(\bar{w} \cdot x_i + \bar{w}_0)$ pe de altă parte: ex. 13.e;
- (P5) 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” ξ_i : ex. 14.a;
- *exemplificarea* noțiunilor de bază: ex. 47, ex. 48 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;
- (P6) 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;
- (P7) o *margine superioară* (engl., upper bound) pentru numărul de *erori* comise la *antrenare* de către C-SVM: ex. 19;

⁹¹⁹ Rezultatul *nu* se menține și în cazul C-SVM: ex. 27.b.

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

- (P8) 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:

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

- (P9) chestiuni legate de *complexitatea computațională* privind clasificatorul C-SVM: ex. 55;
- (•) deducerea *formeii duale* pentru problema SVM cu margine “soft” (C-SVM) de normă \mathcal{L}_2 : ex. 50;
- (•) 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. 51;
- (•) algoritmul SMO (Sequential Minimal Optimization):
deducerea relațiilor de actualizare a variabilelor Lagrange: ex. 22;
exemple de aplicare a algoritmului SMO simplificat: ex. 23, ex. 52;
[SMO pentru *one-class*, *Max Margin SVM*: ex. 31.c;]
- 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. 57.

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. 53, ex. 54;
- *exemplificarea* efectului pe care îl are translatarea datelor în raport cu o axă (Oy) asupra poziției separatorului optimal pentru SVM (în raport cu *funcția-nucleu* folosită): ex. 43;
- 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. 44;
 - (P10) pentru orice set de instanțe distincte și pentru orice etichetare a acestora, există o valoare a hiper-parametrului nucleului RBF (σ) astfel încât SVM obține la antrenare eroare 0: ex. 26. Rezultatul *nu* este valabil și pentru C-SVM;
 - (P11) pentru orice set de instanțe distincte, pentru orice etichetare a acestora și pentru *orice* valoare a hiper-parametrului nucleului RBF (σ), 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. 29;

- avantaje și dezavantaje ale folosirii metodelor de clasificare liniară de tipul SVM, Perceptron etc. și versiunile lor kernel-izate: ex. 56;
- chestiuni recapitulative: ex. 27, ex. 28, ex. 58, ex. 65.

Alte probleme [de optimizare] de tip SVM

- SVM pentru clasificare n -ară (SVM multiclass): ex. 30, ex. 59;
- deducerea *formeii duale* pentru problema *one-class SVM*, versiunea *Max Margin*: ex. 31 (variantea cu margine “hard”), ex. 61 (variantea cu margine “soft”, folosind ν -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. 32 (variantea cu margine “hard”), ex. 62 (variantea cu margine “soft”, folosind ν -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. 32;
- deducerea *formeii duale* pentru problema ν -SVM: ex. 33;
- deducerea *formeii duale* pentru problema SVR (*Support Vector Regression*), folosind funcție de cost / pierdere ε -senzitivă: ex. 34 (cu margine “hard”), ex. 64 (cu margine “soft” și (echivalent) cu funcție de cost ε -senzitivă); exemplificare / aplicare: ex. 63.

6. 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 semiplanele 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. 73, ex. 151 și ex. 152 de la cap. *Fundamente*; ex. 23.

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. 33.⁹²⁰
- (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. 34;
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. 36;
- deducerea regulii de actualizare a ponderilor pentru tipuri particulare de perceptroni: ex. 12, ex. 25.a, ex. 35, 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 semisuma 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.

⁹²⁰Problemele 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;
pentru perceptronul-prag (clasic): ex. 38;
învățare online cu perceptronul-prag de tip Rosenblatt: ex. 37;
- *Perceptronul* kernel-izat [dual]: ex. 19; particularizare pentru cazul nucleului RBF: ex. 39. Perceptronul Rosenblatt, cu termen liber (engl., offset), kernel-izare: ex. 40. Clasificare ternară cu perceptronul Rosenblatt, varianta kernelizată: ex. 41. (Vedeți și ex. 38.c.)

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 sigmoide) cu oricâte niveluri, vedeți p. 101-103;
ex. 20: deducerea regulilor de actualizare a ponderilor în cazul rețelelor cu 2 niveluri, având însă unități cu funcție de activare oarecare (derivabilă);
- aplicare: ex. 21, ex. 43, ex. 44;
- prevenirea overfitting-ului:
folosirea unei componente de tip „moment” în expresia regulilor de actualizare a ponderilor: ex. 46;
regularizare: introducerea unei componente suplimentare în funcția de optimizat: ex. 22;
- cazul folosirii unei funcții de activare de tip tangentă hiperbolică: ex. 45;
- cazul folosirii unei funcții de cost / penalizare / eroare de tip cross-entropie: ex. 48;
- 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. 49.

Rețele neuronale profunde — câteva chestiuni introductive

- analiza convexității unor funcții de cost folosite în învățarea profundă: ex. 53;
- 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. 54.

7. Clusterizare

Sumar

7.0 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 neierarhică).

7.1. Clusterizare ierarhică

7.1.1. 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”:⁹²¹ $d(A, B) = \min\{d(x, y) | x \in A, y \in B\}$

“complete-linkage”:⁹²² $d(A, B) = \max\{d(x, y) | x \in A, y \in B\}$

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

metrica lui Ward: ex. 32, ex. 33, ex. 34.

În general, putem considera $\text{sim}(A, B) = 1/(1 + d(A, B))$ sau chiar $\text{sim}(A, B) = 1/d(A, B)$ dacă lucrăm doar cu 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):

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

⁹²¹Sau: nearest-neighbour.

⁹²²Sau: furthest-neighbour.

7.1.2. Algoritmi de clusterizare ierarhică

- tipuri de algoritmi de clusterizare ierarhică:
bottom-up (*clusterizare aglomerativă*) vs. top-down (*clusterizare divizivă*);
- pseudo-coduri (cf. Manning & Schütze, *Foundations of Statistical Natural Language Processing*, 2002, pag. 502):

bottom-up:

```

Given: a set  $X = \{x_1, \dots, x_n\}$  of objects
          a function sim:  $\mathcal{P}(X) \times \mathcal{P}(X) \rightarrow R$ 
for  $i = 1, n$  do
     $c_i = \{x_i\}$  end
 $C = \{c_1, \dots, c_n\}$ 
 $j = n + 1$ 
while  $|C| > 1$ 
     $(c_{n_1}, c_{n_2}) = \operatorname{argmax}_{(c_u, c_v) \in C \times C} \mathbf{sim}(c_u, c_v)$ 
     $c_j = c_{n_1} \cup c_{n_2}$ 
     $C = C \setminus \{c_{n_1}, c_{n_2}\} \cup \{c_j\}$ 
     $j = j + 1$ 

```

top-down:

```

Given: a set  $X = \{x_1, \dots, x_n\}$  of objects
          a function coh:  $\mathcal{P}(X) \rightarrow R$ 
          a function split:  $\mathcal{P}(X) \rightarrow \mathcal{P}(X) \times \mathcal{P}(X)$ 
 $C = \{X\} (= \{c_1\})$ 
 $j = 1$ 
while  $\exists c_i \in C$  such that  $|c_i| > 1$ 
     $c_u = \operatorname{argmin}_{c_v \in C} \mathbf{coh}(c_v)$ 
     $c_{j+1} \cup c_{j+2} = \mathbf{split}(c_u)$ 
     $C = C \setminus \{c_u\} \cup \{c_{j+1}, c_{j+2}\}$ 
     $j = j + 2$ 

```

- 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. 27-31 (pentru bottom-up), respectiv ex. 6 (pentru top-down);
- implementări: ex. 34, ex. 36, ex. 37.

7.1.3 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. 30;
- (P1) dacă atunci când folosim “single-linkage” și “complete-linkage” se obțin dendrograme identice, *nu* rezultă în mod neapărat că folosind “average-linkage” vom obține aceeași dendrogramă: ex. 3.b;
- (P2) 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;

- (P3) 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);
- (P4) algoritmul de clusterizare aglomerativă la al cărui pseudo-cod am făcut referire mai sus are complexitate $\mathcal{O}(n^3)$: ex. 27; 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”:

(P5) 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 a] unui cluster nou format, precum și calculul „distanței” dintre două cluster se pot face în timp constant: ex. 35.

7.2. Clusterizare partițională

7.2.1 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, μ este K -configurație de centroizi, iar $\mu_{C(x_i)}$ este centroidul cel mai apropiat de x_i : ex. 12.

7.2.2 Algoritmul K -means

- pseudo-cod (cf. Manning & Schütze, *op. cit.*, pag. 516):

Given: a set $X = \{x_1, \dots, x_n\} \subseteq \mathcal{R}^m$,
 a distance measure d on \mathcal{R}^m ,
 a function for computing the mean $\mu : \mathcal{P}(\mathcal{R}^m) \rightarrow \mathcal{R}^m$

Select (arbitrarily) k initial centers f_1, \dots, f_k in \mathcal{R}^m ;
 while the *stopping criterion* is not satisfied
 for all clusters c_j do $c_j = \{x_i \mid \forall f_l \ d(x_i, f_j) \leq d(x_i, f_l)\}$ end
 for all means f_j do $f_j \leftarrow \mu(c_j)$ end

alternativ, vedeți enunțul ex. 12 (sau, echivalent, folosind variabile-indicator: ex. 43);
- exemple de aplicare: ex. 7-11, ex. 17.a, ex.21.a, ex. 22.a, ex. 38, ex. 39.
- exemple de *euristici pentru inițializarea centroizilor*: 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. 47.

- 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*:
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 centroizilor: 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. 46;
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. 44.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 liniare 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. 48;
 (P4) rezultatele lui *K*-means pot fi influențate de prezența outlier-elor: ex. 10; în astfel de cazuri este de preferat să se folosească distanța Manhattan în locul distanței euclidiene: ex. 46.b.
- o euristică pentru alegerea unei valori convenabile / „naturale” pentru K (pentru un dataset dat) — criteriul “elbow”: ex. 42 (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. 46;
- implementare: ex. 51.

7.2.3 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. 42);
 - (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. 44.
- 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,

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

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 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. 43).

7.3. Clusterizare prin modelare probabilistă

7.3.1 Noțiuni preliminare

- variabile aleatoare (discrete, resp. continue);
media, varianța și covarianța 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. 20 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 multidimensionale: vedeți secțiunea corespunzătoare din *Sumarul capitolului de Fundamente*;
- 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. 26 și ex. 104 de la capitolul de *Fundamente*;
 - o mixtură de distribuții gaussiene multidimensionale: ex. 108 de la capitolul de *Fundamente*;
 - o mixtură de distribuții oarecare: ex. 109 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 θ se notează parametrii respectivei distribuții. Exemplificare: ex. 40.abd, ex. 39 de la capitolul de *Fundamente*;
- MLE (Maximum Likelihood Estimation): estimarea [valorilor] parametrilor unei distribuții probabiliste în sensul maximizării verosimilității datelor disponibile. Exemplificare: capitolul de *Fundamente*, ex. 40-49, ex. 114-125. Aplicare în cazul distribuției gaussiene unidimensionale: ex. 15.ab de la capitolul *Clasificare bayesiană*;
- regula de decizie pentru algoritmul Bayes [Naiv] Gaussian:⁹²³ pentru cazul unidimensional: ex. 15, ex. 39, ex. 16, și ex. 40 de la capitolul *Clasificare bayesiană*;
- 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. 65.

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

- pseudo-cod:⁹²⁴
 - cazul unidimensional, varianta când doar parametrul μ este lăsat liber: ex. 15 (cf. *Machine Learning*, Tom Mitchell, 1997, pag. 193); aplicare: ex. 16, ex. 17.b, ex. 52, ex. 53;
 - cazul unidimensional, varianta când toți parametrii (π , μ și σ) sunt lăsați liberi: ex. 18; aplicare: ex. 17.c, ex. 54, ex. 55;
 - cazul unidimensional, alte variante: [ex. 19] ex. 56, ex. 57, ex. 58;

⁹²³În cazul separării lineare, în literatura de specialitate se folosește termenul de *analiză discriminativă gaussiană*.

⁹²⁴Nu doar pentru pseudo-cod, ci și (sau, mai ales) pentru o privire de ansamblu unitară, atât pentru cazul unidimensional cât și pentru cazul multidimensional (ambele urmând a fi sistematizate mai jos), puteți consulta documentul *An Introduction to Expectation-Maximization*, de Dahua Lin, MIT, ML course 6768, 2012 fall.

cazul multidimensional, varianta când toți parametrii (π , μ și Σ) sunt lăsați liberi: ex. 24;
 cazul multidimensional, alte variante: ex. 20, ex. 60;
 aplicarea algoritmului EM/GMM, cazul bidimensional: ex. 21.b, ex. 22.b, ex. 23, ex. 25, ex. 61, ex. 62, ex. 63;

- schema algoritmică EM: vedeți Tom Mitchell, *Machine Learning* book, 1997, pag. 194-195; ex. 19;
- ca algoritm de *învățare statistică*:
 algoritmul EM poate fi văzut ca o metodă de estimare a parametrilor (engl., parameter fitting);
- ca algoritm *per se*:
 - *algoritm iterativ*: pleacă de la o soluție (instantiere pentru parametri) aleasă eventual în mod arbitrar / aleatoriu și o „îmbunătățește“ la fiecare iterație;
 - *algoritm de căutare*: se caută într-o anumită clasă de modele probablistice (parametrizate) un model care să satisfacă *principiul verosimilității maxime*;⁹²⁵
 - *algoritm 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 ;

mai precis, 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 (θ), se alege în mod arbitrar, iar apoi $\theta^{(t+1)} = \arg\max_{\theta} Q_t(\theta|\theta^{(t)})$;

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*;

⁹²⁵ Acesta este cazul general. În cazuri particulare, acest principiu poate fi înlocuit cu *principiul probabilității maxime a posteriori*. Vedeți problema 3 de la capitolul *Algoritmul EM*.

învățare în prezența unor variabile aleatoare neobservabile(!);
[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. 17.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. 19.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. 17.c, pag. 801, sau la ex. 63.c și ex. 64.c).

- Comparativ cu algoritmul K -means,
 - (P5) algoritmul EM/GMM este în general mai lent — mișcarea centroizilor poate explora într-o manieră mai fină spațiul (vedeți de exemplu ex. 21) —, dar din acest motiv el poate să obțină uneori rezultate mai bune / convenabile (vedeți spre exemplu ex. 22);
 - (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. 17.b);
 - (P7) EM/GMM este mai robust la influența outlier-elor.

7.3.3 Alte proprietăți ale algoritmului EM/GMM

- (P8) Legătura dintre algoritmul K -means și algoritmul EM/GMM:⁹²⁶ 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. 60;
- 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. 65;
- Algoritmului EM/GMM *semi-supervizat*: ex. 66;
o variantă a algoritmului EM/GMM semi-supervizată pentru cazul când matricele de covarianță sunt de forma $\sigma_k^2 I$ și este satisfăcută presupuziția de independență condițională de tip Bayes Naiv: ex. 67 (are loc o legătură interesantă cu clasificatorul Bayes Naiv gaussian).

⁹²⁶Formularea se referă la cazul multidimensional, dar proprietatea este valabilă și în cazul unidimensional.

8. Schema algoritmică EM

Sumar

Noțiuni preliminare

- distribuții probabiliste: vedeți secțiunea *Distribuții probabiliste uzuale* de la capitolul de *Fundamente*;
- estimarea parametrilor distribuțiilor probabiliste în sensul verosimilității maxime (MLE) și respectiv în sensul probabilității maxime a posteriori (MAP): vedeți secțiunea *Estimarea parametrilor unor distribuții probabiliste* de la capitolul de *Fundamente*;
- mixturi de distribuții probabiliste: vedeți ex. 26 și ex. 104 de la capitolul de *Fundamente*, ex. 5.AB și ex. 23 de la prezentul capitol;⁹²⁷
- metoda maximizării alternante pe coordonate (engl., *coordinate ascent*) pentru rezolvarea unor probleme de optimizare: ex. 1;⁹²⁸
- *metoda multiplicatorilor lui Lagrange* pentru rezolvarea unor probleme de optimizare cu restricții: vedeți secțiunea *Metode de optimizare în învățarea automată* de la capitolul de *Fundamente*, precum și ex. 9, ex. 10, ex. 12, ex. 24, ex. 25 și ex. 32 de la prezentul capitol.

Schema algoritmică EM

- pseudo-cod: *Machine Learning*, Tom Mitchell, 1997, pag. 194-195;
- fundamentare teoretică:
 - ex. 1: pentru funcția de log-verosimilitate a datelor complete, există o margine inferioară, $F(q, \theta)$; algoritmul EM face maximizarea acestei margini inferioare aplicând metoda [iterativă a] creșterii alternative pe coordonate (engl., *coordinate ascent*);
 - ex. 2: monotonia valorilor funcției de log-verosimilitate a datelor complete, care sunt calculate la iterații succesive ale lui EM;⁹²⁹
 - ex. 3: MAP EM – algoritmul EM pentru *estimarea* nu în sens MLE (cum este cazul adeseori), ci *în sens MAP*; pentru exemplificare, vedeți ex. 29.B;
 - ex. 20: algoritmul EM semi-supervizat (particularizare pentru cazul mixturilor de distribuții Bernoulli: ex. 35);⁹³⁰
 - ex. 21: “hard” EM – algoritmul EM cu asignare “hard” a instanțelor la clustere;
 - ex. 22: algoritmul EM generalizat (engl., Generalized EM, GEM).

EM pentru modelarea de mixturi de distribuții probabiliste

⁹²⁷Pentru mixturi de distribuții gaussiene vedeți secțiunea *Algoritmul EM pentru modele de mixturi gaussiene* de la capitolul de *Clusterizare*.

⁹²⁸Vedeți, de asemenea, utilizarea aceleiași metode de optimizare (eventual pentru minimizare în loc de maximizare) în cazul altor algoritmi de învățare automată: pentru algoritmul AdaBoost, ex. 22, ex. 28 și ex. 29 de la capitolul de *Arbori de decizie*; pentru algoritmul K-means, ex. 12 de la capitolul de *Clusterizare*; în sfârșit, pentru algoritmul SMO, ex. 22 de la capitolul de *Mașini cu vectori-suport*.

⁹²⁹Acestea, plus alte câteva *proprietăți generale* ale schemei algoritmice EM au fost deja sintetizate sub forma proprietăților (P0)-(P3) din secțiunea *Algoritmul EM pentru clusterizare* de la *Sumarul* capitolului de *Clusterizare*.

⁹³⁰Pentru cazul mixturilor de distribuții gaussiene, vedeți ex. 66 de la capitolul *Clusterizare*.

- mixturi de distribuții *Bernoulli*: ex. 5, ex. 6 și ex. 23;
mixturi de *vectori* de distribuții *Bernoulli*, cu presupunerea de independență condițională a atributelor de intrare în raport cu atributul de ieșire (eticheta): ex. 7, ex. 8;
o versiune particulară — clusterizare în interiorul claselor —, cu *aplicare* la recunoașterea cifrelor scrise de mână: ex. 26;
- mixturi de distribuții *categoriale*: ex. 9 și ex. 24;
aplicații: [EM pentru] *dezambiguizarea semantică* a cuvintelor dintr-un document (engl., word sense disambiguation): ex. 10 și respectiv pentru clusterizare de documente (engl., topic model): ex. 27;
- mixturi de *vectori* de distribuții *categoriale*, cu presupunerea de independență condițională a atributelor de intrare în raport cu atributul de ieșire (eticheta): ex. 25 (algoritmul *Bayes Naiv nesupervizat*, cu *asignare "soft"* a instanțelor la clustere);
- mixturi de [vectori de] distribuții *Poisson*: ex. 31, ex. 32;
- mixturi de distribuții *Gamma*: ex. 33;
- EM pentru *estimarea probabilității de selecție* a unei componente din cadrul unei mixturi (i.e., combinație liniară) de două distribuții probabiliste oarecare: ex. 16.

EM pentru estimarea parametrilor unor distribuții probabiliste

- EM pentru estimarea parametrilor unor distribuții *binomiale*: ex. 28;⁹³¹
- EM pentru estimarea parametrilor unor distribuții *multinomiale* (care se definesc cu ajutorul uneia sau mai multor distribuții *categoriale*): ex. 11, ex. 12, ex. 29, ex. 30.

EM pentru estimarea parametrilor unei distribuții, atunci când o parte din date lipesc

- cazul distribuției *Poisson*: ex. 15.

EM pentru estimarea parametrilor unei *sume* de două distribuții⁹³²

- cazul distribuțiilor *exponențiale*: ex. 13;
- cazul distribuțiilor *gaussiene*: ex. 14.

⁹³¹În acest exercițiu, apar trei distribuții binomiale, dintre care una se definește cu ajutorul unei mixturi de două distribuții Bernoulli.

⁹³²Explicație: ne referim aici la aplicarea algoritmului EM pentru estimarea parametrilor a două distribuții probabiliste atunci când se dau instanțe care sunt generate prin însumarea unor perechi de valori generate de cele două distribuții;

Alte instanțe ale schemei algoritmice EM

- algoritmul EM pentru învățare supervizată; cazul mixturilor de regresori liniari: ex. 34.

Alte probleme

- chestiuni metodologice (relativ la inițializarea parametrilor): ex. 36;
- probleme recapitulative (A/F): ex. 4, ex. 17, ex. 18, ex. 19 și ex. 37.

9. Modele Markov ascunse

Sumar

- Noțiuni preliminare: programare dinamică, schema algoritmică EM.
- Verificarea înțelegerii unor noțiuni de bază în referitoare la modelele Markov ascunse (engl., Hidden Markov Models, HMM): ex. 1-4, ex. 15-17.
- Model Markov vizibil: exemplificare, legătura cu HMM: ex. 5.
- HMM ca model probabilist total: ex. 14.
- Algoritmul Forward:⁹³³
exemple de aplicare: ex. 6, ex. 7, ex. 9.a;
demonstrarea formulei de la pasul inductiv: ex. 18.a;
calcularea probabilității de emisie a unei secvențe, folosind probabilitățile Forward: ex. 18.b.
- Algoritmul Backward:⁹³⁴
demonstrarea formulei de la pasul inductiv: ex. 8;
exemplu de aplicare: ex. 9.a.
- Algoritmul Viterbi:⁹³⁵
exemplu de aplicare: ex. 9.b;
determinarea căii celei mai probabile de generare a unei secvențe: ex. 9.c;
determinarea probabilității ca un anumit simbol [din secvența de semnale] să fi fost generat într-o anumită stare: ex. 9.d;
o variație pe tema algoritmului Viterbi: ex. 20.
- Algoritmul Forward-Backward / Baum-Welch / EM pentru HMM:⁹³⁶
exemplu de aplicare: ex. 11.b;
demonstrarea formulei necesare pentru calcularea mediilor variabilelor neobservabile care corespund tranzițiilor: ex 10;
demonstrarea faptului că algoritmul Forward-Backward lasă neschimbate probabilitățile de tranziție sau de emisie care sunt nule [la inițializare]: ex 12.
- Demonstrarea unei formule alternative — în raport cu formulele bazate pe probabilitățile Forward și respectiv probabilitățile Backward — pentru calcularea probabilității de emisie a unei secvențe: ex. 19.
- HMM cu emisii gaussiene: ex. 13, ex. 22.

⁹³³Pentru calculul probabilităților Forward, notate cu $\alpha_i(t) = P(O_1, \dots, O_t, X_{t+1} = S_i)$.

⁹³⁴Pentru calculul probabilităților Backward, notate cu $\beta_i(t) = P(O_t O_{t+1} \dots O_T | X_t = S_i)$.

⁹³⁵Pentru calculul cantităților $\delta_i(t) = \max_{X_1 \dots X_{t-1}} P(X_1 \dots X_{t-1}, O_1 \dots O_{t-1}, X_t = s_i)$.

⁹³⁶Pentru „învățarea” parametrilor (i.e., a probabilităților de tranziție și respectiv de emisie) ai / ale unui HMM.