

Universitatea „Alexandru Ioan Cuza” Iași

Facultatea de Informatică



LUCRARE DE LICENȚĂ

**O analiză asupra algoritmilor de
clasificare discriminativi și generativi și
aplicarea acestora în diagnoza cancerului
la sân**

propusă de

Giurgilă Andreea

Sesiunea: Iulie, 2018

Coordonator științific

Conferențiar, Dr. Liviu Ciortuz

Universitatea „Alexandru Ioan Cuza” Iași
Facultatea de Informatică

**O analiză asupra algoritmilor de
clasificare discriminativi și generativi și
aplicarea acestora în diagnoza cancerului
la sân**

Giurgilă Andreea

Sesiunea: Iulie, 2018

Coordonator științific
Conferențiar, Dr. Liviu Ciortuz

Avizat,
Îndrumător Lucrare de Licență

Titlul, Numele și prenumele _____

Data _____ Semnătura _____

**DECLARAȚIE privind originalitatea conținutului lucrării de
licență/diplomă/disertație/absolvire**

Subsemnata **GIURGILĂ ANDREEA** domiciliul în **ADJUDENI, NEAMT**

născut(ă) la data de **03 FEB 1997**, identificat prin CNP **2970203270831**, absolvent(a) al(a) Universității „Alexandru Ioan Cuza” din Iași, Facultatea de Informatică promoția 2018, declar pe propria răspundere, cunoscând consecințele falsului în declarații în sensul art. 326 din Noul Cod Penal și dispozițiile Legii Educației Naționale nr. 1/2011 art.143 al. 4 și 5 referitoare la plagiat, că lucrarea de licență cu titlul *“O analiză asupra algoritmilor de clasificare discriminativi și generativi și aplicarea acestora în diagnoza cancerului la sân”* elaborată sub îndrumarea dl. Conf. Dr. Liviu Ciortuz, pe care urmează să o susțină în fața comisiei este originală, îmi aparține și îmi asum conținutul său în întregime.

De asemenea, declar că sunt de acord ca lucrarea mea de licență/diplomă/disertație/absolvire să fie verificată prin orice modalitate legală pentru confirmarea originalității, consimțind inclusiv la introducerea conținutului său într-o bază de date în acest scop.

Am luat la cunoștință despre faptul că este interzisă comercializarea de lucrări științifice în vederea facilitării falsificării de către cumpărător a calității de autor al unei lucrări de licență, de diploma sau de disertație și în acest sens, declar pe proprie răspundere că lucrarea de față nu a fost copiată ci reprezintă rodul cercetării pe care am întreprins-o.

Data azi, Semnătură student

DECLARAȚIE DE CONSIMȚĂMÂNT

Prin prezenta declar că sunt de acord ca Lucrarea de licență cu titlul „*O analiză asupra algoritmilor discriminativi și generativi și aplicarea acestora în diagnoza cancerului la sân*”, codul sursă al programelor și celelalte conținuturi (grafice, multimedia, date de test etc.) care însoțesc această lucrare să fie utilizate în cadrul Facultății de Informatică.

De asemenea, sunt de acord ca Facultatea de Informatică de la Universitatea Alexandru Ioan Cuza Iași să utilizeze, modifice, reproducă și să distribuie în scopuri necomerciale programele calculator, format executabil și sursă, realizate de mine în cadrul prezentei lucrări de licență.

Iași, 29.06.2018

Giurgilă Andreea

(semnătura în original)

MULȚUMIRI

În primul rând, aș dori să mulțumesc coordonatorului meu științific, Dr. Conf. Liviu Ciortuz. Fără ajutorul și implicarea lui nu aș fi reușit să scriu această lucrare. De asemenea aș dori să-i mulțumesc pentru introducerea pe care mi-a oferit-o în acest domeniu în cadrul cursului de Învățare Automată de la Facultatea de Informatica din Iași. Acest curs a fost cel care m-a determinat să aprofundez noțiunile introduse la curs și să le extind într-o lucrare de cercetare.

Un ajutor foarte mare în scrierea acestei lucrări au fost cursurile și materialele domnilor profesori Andrew Ng. (Stanford) și Tom Mitchell (Carnegie Mellon University). Deși aceștia nu au fost profesorii mei, resursele postate online m-au ajutat în tratarea anumitor subiecte.

De asemenea, sunt recunoscătoare Facultății de Informatică din Iasi pentru oferirea resurselor necesare scrierii acestei lucrari.

Cuprins

1	Introducere	8
2	Modelul Generativ	9
2.1	Introducere capitol	9
2.2	Teorema lui Bayes pentru învățarea de etichete (ipoteze)	9
2.2.1	Estimarea probabilităților (MLE și MAP)	9
	Estimarea MLE	10
	Estimarea MAP	11
2.3	Algoritmul Bayes Naiv	13
2.3.1	Noțiunea de Independența Condițională	13
2.3.2	Introducere Bayes Naiv	13
2.3.3	Relația de echivalență în exprimare dintre Bayes Naiv și Regresia Logistica pentru variabile discrete	14
2.3.4	Algoritmul Bayes Naiv pentru variabile continue	15
2.3.5	Bayes Naiv Gaussian	16
2.4	Concluzie capitol	17
3	Modelul Discriminativ	18
3.1	Introducere capitol	18
3.2	Regresia Logistica	18
3.2.1	Introducere Regresie Logistica	18
3.2.2	Estimarea parametrilor	19
3.2.3	Actualizarea valorilor w_i	20
3.2.4	Teorema de convergență pentru Regresia Logistica	20
3.2.5	Regularizarea regresiei logistice	21
3.2.6	Regresia Logistica Kernelizată	22
3.3	Mașini cu Vector Suport	24
3.4	Concluzie capitol	25

4	Exemple Aplicative	27
4.1	Aplicare Regresie Logistica	27
4.1.1	Vizualizarea graniței de decizie	28
4.1.2	Regresia Logistica MultiClass	29
4.1.3	Alte Observatii	29
4.2	Aplicare Bayes Naiv vs Regresia Logistica	30
4.3	Concluzii Capitol	31
5	Prezicerea Cancerului la san	32
5.1	Introducere	32
5.2	Date Utilizate	33
5.2.1	Alegerea Atributelor	33
5.3	Metode și Rezultate	35
5.4	Concluzii	37
5.5	Direcții Viitoare	38
6	Bibliografie	39

1 Introducere

Odată cu creșterea exponențială a cantității de date existente în lume, necesitatea unor metode potrivite de analiză a acestor date crește și ea. Prin analiză înțelegem învățarea unui model dintr-un set de date cu ajutorul căruia să putem prezice pentru un set de date noi un anumit comportament. Prin crearea unor modele probabiliste, învățarea Automată are ca scop automatizarea acestui proces de analiză.

Lucrarea de față este structurată pe patru capitole care vor ajuta la înțelegerea modului în care algoritmi generativi diferă față de cei discriminativi și cum se pot aplica aceștia în rezolvarea unei probleme din viața reală și anume diagnoza cancerului la sân.

Capitolul I cuprinde o scurtă descriere a modului în care lucrează un algoritm de clasificare generativ. În special, voi prezenta algoritmul Bayes Naiv și o variație a acestuia, Bayes Naiv Gaussian. De asemenea voi arăta cum forma unuia dintre parametrii acestui algoritm se aseamănă cu un parametru dintr-un algoritm discriminativ - Regresia Logistică.

Capitolul II prezintă modul de lucru al algoritmilor de clasificare discriminativi. În special, acest capitol prezintă doi algoritmi: Regresia Logistică și Mașinile cu Vector Suport.

Capitolul III descrie câteva aplicații ale algoritmilor specificați mai sus, punând în evidență diferența dintre aceștia. Voi aplica un algoritm de clasificare discriminativ și unul generativ pe același set de date și voi compara rezultatele.

În **Capitolul IV** voi aplica algoritmi discutați în această lucrare în rezolvarea unei probleme din viața reală: diagnoză a cancerului la sân. Voi prezenta modul de abordare al acestei probleme, cum se face selecția de atribute și voi interpreta rezultatele obținute.

2 Modelul Generativ

2.1 Introducere capitol

Un model generativ este un model care, folosindu-se de un model probabilistic, descrie modul în care sunt generate datele.

Presupunem că avem o problemă de clasificare în care ne dorim să deosebim între motociclete și mașini. Pentru a învăța pe cineva să recunoască o mașină desenăm pe o foaie de hârtie modelul unei mașini și pe o foaie separată cel al unei motociclete. Când ne confruntăm cu problemă de a clasifica un nou obiect, trebuie doar să decidem cu care dintre modelele desenate anterior este mai asemănător obiectul. Acest tip de clasificare se încadrează în categoria algoritmilor de învățare generativi. Mai formal, numim algoritm generativ un algoritm care învață o funcție scor $f(x, y)$ a datelor de intrare x și etichetele y . Pentru a prezice o nouă instanță y , algoritmul alege y_i pentru care funcția $f(x, y_i)$ este maximă.

2.2 Teorema lui Bayes pentru învățarea de etichete (ipoteze)

Fiind dată o instanță x_i și eticheta corespunzătoare x_i , formula lui Bayes ne prezintă o relație între probabilitatea lui x_i atunci când s-a observat deja eticheta y_i (probabilitatea lui x_i a posteriorii observării y_i) și probabilitatea ca rezultatul clasificării să fie y_i , ținând cont ca instanța de clasificat este x_i (probabilitatea lui y_i a posteriorii lui x_i). Mai precis:

$$p(x_i | y_i) = \frac{p(y_i | x_i)p(x_i)}{p(y_i)} \quad (1)$$

Unde $p(x_i)$ este probabilitatea a priori a lui x_i iar $p(y_i)$ este probabilitatea a priori a lui $p(x_i)$.

2.2.1 Estimarea probabilităților (MLE și MAP)

Atunci când ne dorim să estimăm un model de distribuție probabilistic al datelor de antrenare, există două metode populare de abordare: estimarea în sensul verosimilității maxime (MLE - en: Maximum Likelihood Estimation) și estimarea în sensul probabilității maxime a posteriorii (MAP

- en: Maximum A Posterior).

Pentru a deosebi între cele două moduri de abordare ne vom folosi de un exemplu simplu*: aruncarea unei monezi. Vom nota cu M variabila care ne dă rezultatul aruncării, cu $p = p(M=\text{ban})$, $p \in (0,1)$, probabilitatea ca rezultatul aruncării să fie ban ($1-p=p(M=\text{stema})$ va fi probabilitatea ca rezultatul să fie stema) și cu \hat{p} probabilitatea estimată ca rezultatul să fie ban.

Dupa n aruncări, obținem de a_1 ori ban și de a_2 ori stema. Intuitiv, deducem ca $p = \frac{a_1}{a_1+a_2}$. Acest mod de gândire este unul bun atunci când datele de antrenare sunt suficiente. În caz contrar, acest algoritm ne poate conduce spre rezultate eronate. De exemplu, dacă aruncând moneda de 4 ori obținem în mod consecutiv ban, vom deduce ca $p(M=\text{ban})=1$ ceea ce este puțin probabil să se întâmple dacă stim că folosim o monedă cinstită care are $P(M=\text{ban}) \sim 0.5$. Acest argument stă la baza celei de-a doua metode de abordare.

A doua metodă ne permite să ne folosim de presupuneri sau cunoștințe a priori despre monedă pentru a obține un rezultat cât mai aproape de adevăr. În cazul nostru, putem introduce un număr imaginar b_1 de rezultate ban respectiv un număr b_2 de aruncări în urma cărora a fost observată stema. Aceste numere vor fi alese conform presupunerii a priori pe care o avem despre monedă. De exemplu, presupunând că moneda este una cinstită și că $P(M=\text{ban})=0.5$ putem alege $b_1 = b_2 = 5$. Astfel noua estimare a lui \hat{p} devine: $\hat{p} = \frac{a_1+b_1}{(a_1+b_1)+(a_2+b_2)}$.

În continuare, vom demonstra că cele două metode prezentate mai sus corespund estimării probabilității în sens MLE, respectiv MAP.

2.2.1.1 Estimarea MLE Estimarea MLE calculează parametrii unui model care explică cel mai bine datele de antrenare D . Prin definiție:

$$\hat{p} = \operatorname{argmax}_p P(D | p) \quad (2)$$

Definim variabila aleatoare M astfel:

$$M = \begin{cases} 1, & \text{daca rezultatul dupa o aruncare e ban} \\ 0, & \text{altfel} \end{cases}$$

De asemenea, notam cu θ probabilitatea adevărată ca rezultatul după o aruncare este ban. Observăm că, dacă avem o singură aruncare, $p(D | \theta) = \theta$ dacă $M=1$ respectiv $p(D | \theta)=1-\theta$ dacă $M=0$. Presupunem că după n aruncări consecutive și independente obținem de b_1 ori rezultatul $M = 1$ și de b_2 ori rezultatul $M = 0$. Putem calcula $P(D | \theta)$ înmulțind fiecare probabilitate. Vom avea:

$$P(D | \theta) = \theta^{b_1} (1 - \theta)^{b_2} \quad (3)$$

Folosind estimatorul MLE, ne dorim să determinăm valoarea lui care maximizează $P(D | \theta)$. Acest lucru este echivalent cu maximizarea logaritmului:

$$f(\theta) = \operatorname{argmax}_{\theta} P(D | \theta) = \operatorname{argmax}_{\theta} \ln P(D | \theta) \quad (4)$$

Pentru a găsi maximum, calculăm derivata funcției de mai sus și o egalăm cu 0.

$$\begin{aligned} \frac{d}{d\theta} f(\theta) &= \frac{d}{d\theta} \ln P(D | \theta) \\ &= \frac{d}{d\theta} \ln(\theta^{b_1} (1 - \theta)^{b_2}) \\ &= \frac{d}{d\theta} [\ln \theta^{b_1} + \ln(1 - \theta)^{b_2}] \\ &= \frac{d}{d\theta} [b_1 \ln \theta + b_2 \ln(1 - \theta)] \\ &= b_1 \frac{d}{d\theta} \ln \theta + b_2 \frac{d}{d\theta} \ln(1 - \theta) \\ &= b_1 \frac{1}{\theta} + b_2 \frac{-1}{1 - \theta} \end{aligned}$$

Nota: Mai sus am folosit derivatele:

$$\begin{aligned} \frac{d}{dx} \ln x &= \frac{1}{x} \\ \frac{d}{dx} \ln g(x) &= \frac{1}{g(x)} \frac{d}{dx} g(x) \end{aligned}$$

Observăm că am obținut același rezultat ca mai sus deci algoritmul intuitiv de mai sus este echivalent cu estimarea în sens MLE.

2.2.1.2 Estimarea MAP

Estimatorul MAP ne permite sa folosim cunoștințe anterioare despre distribuția a priori $p()$ (folosim aceeași notație ca mai sus). Aceste cunoștințe pot rezulta din experienta noastră, din domeniul de specialitate al problemei sau pot fi intuitive. Astfel, estimarea MAP se va folosi de probabilitatea a posteriorii:

$$\hat{p} = \operatorname{argmax}_p P(D | p) \quad (5)$$

Pentru a ajunge la forma finala la formula de mai sus am folosit formula lui Bayes, astfel:

$$P(p | D) = \frac{P(D | p)P(p)}{P(D)} \quad (6)$$

Numitorul a fost mai apoi ignorat întrucât nu depinde de p . In general, atunci când ne confruntam cu rezultate consecutive ale unei variabile aleatoare Bernoulli, datele urmează o distribuție Beta. Prin urmare putem scrie $P(p)$ astfel:

$$P(p) = \operatorname{Beta}(b_1, b_2) = \frac{p^{b_1-1}(1-p)^{b_2-1}}{B(b_1, b_2)} \quad (7)$$

Unde $B(b_1, b_2)$ este un factor de normalizare pe care il vom ignora in continuare întrucât nu depinde de p . Cu aceasta noua informație putem sa rescriem prima ecuație astfel:

$$\begin{aligned} p^{MAP} &= \operatorname{argmax}_p P(D | p)P(p) \\ &= \operatorname{argmax}_p p^{a_1}(1-p)^{a_0} \frac{p^{b_1-1}(1-p)^{b_2-1}}{B(b_1, b_2)} \\ &= \operatorname{argmax}_p \frac{p^{a_1+b_1-1}(1-p)^{a_0+b_2-1}}{B(b_1, b_2)} \\ &= \operatorname{argmax}_p p^{a_1+b_1-1}(1-p)^{a_0+b_2-1} \end{aligned}$$

Observam ca pentru a afla valoarea p care maximizează funcția de mai sus ne putem folosi de rezultatul obținut la Estimarea MLE, înlocuind $b_1 = a_1 + b_1 - 1, b_2 = a_2 + b_2 - 1$. In final obținem:

$$p^{MAP} = \operatorname{argmax}_p P(D | p)P(p) = \frac{(a_1 + b_1 - 1)}{(a_1 + b_1 - 1) + (a_0 + b_2 - 1)} \quad (8)$$

,același rezultat pe care l-am obținut și în algoritmul intuitiv inițial dacă facem substituția $b_1 = \beta_1 - 1$ și $b_2 = \beta_2 - 1$. Deci, folosind estimarea în sens MAP - mai exact dorind să găsim valoarea lui $P(\theta)$ care maximizează $P(\theta | D)$ - am obținut același rezultat

2.3 Algoritmul Bayes Naiv

2.3.1 Noțiunea de Independența Condițională

Spunem că două variabile aleatoare X, Y sunt independente condițional în raport cu o a treia variabilă aleatoare Z dacă următoarea egalitate este îndeplinită:

$$P(X, Y | Z) = P(X | Z)P(Y | Z) \quad (9)$$

2.3.2 Introducere Bayes Naiv

Fiind dat un set de date de antrenare $X = X_1, X_2, X_3, \dots, X_n$ și etichetele corespunzătoare reprezentarea de variabilă discretă $Y = Y_1, Y_2, \dots, Y_n$, algoritmul Bayes Naiv de folosește de formula lui Bayes enunțată la 2.1.1 și de presupunerea de independență condițională a datelor de antrenare pentru a estima eticheta Y_k a unei noi instanțe $X_k = x_1, x_2, \dots$. Pentru a face acest lucru, algoritmul va învăța o distribuție asupra probabilităților pentru fiecare valoare posibilă a etichetei Y_k în raport cu noua instanță de clasificat $X_k = x_1, x_2, \dots$ și în final va alege cea mai probabilă etichetă.

Conform 2.2.1, putem scrie probabilitatea ca Y_k să aibă valoarea y_k pentru instanța $X_k = x_1, x_2, x_3$ astfel:

$$P(Y_k = y_k | x_1, x_2, \dots) = \frac{P(Y_k = y_k)P(x_1, x_2, x_3, \dots | Y = y_k)}{\sum_j P(Y_k = y_j)P(x_1, x_2, x_3, \dots | Y = y_j)} \quad (10)$$

Nota: Pentru exprimarea numitorului $P(x_1, x_2, x_3, \dots)$ s-a folosit Formula Probabilității Totale

Presupunând că x_1, x_2, x_3 sunt independente condițional în raport cu y_k vom rescrie ecuația de mai sus astfel: $P(Y_k = y_k | x_1, x_2, \dots) = \frac{P(Y_k = y_k) \prod_i P(x_i | Y = y_k)}{\sum_j P(Y_k = y_j) \prod_i P(x_i | Y = y_j)}$

În final, pentru a afla eticheta cea mai probabilă trebuie doar să aflăm maximum dintre toate valorile posibile pentru Y :

$$Y_k = \operatorname{argmax}_k P(Y_k = y_k | x_1, x_2, \dots) = \operatorname{argmax}_k \frac{P(Y_k = y_k) \prod_i P(x_i | Y = y_k)}{\sum_j P(Y_k = y_j) \prod_i P(x_i | Y = y_j)} \quad (11)$$

întrucât numitorul nu depinde de y_k îl putem ignora:

$$Y_k = \operatorname{argmax}_k P(Y_k = y_k) \prod_i P(x_i | Y = y_k) \quad (12)$$

2.3.3 Relația de echivalență în exprimare dintre Bayes Naiv și Regresia Logistica pentru variabile discrete

În continuare datele de intrare vor avea aceeași formă descrisă mai sus. În plus, vom presupune că variabilele X_i sunt boolene. Vom demonstra că distribuția condițională $P(Y | X)$ are forma funcției logistice de parametru w , descrisă mai în detaliu în secțiunea 3, și anume:

$$P(Y = 1 | X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)} \quad (13)$$

Vom introduce următoarele notații: $P(X_i = 1 | Y = 1) = \theta_{i1} \Rightarrow P(X_i = 0 | Y = 1) = 1 - \theta_{i1}$, $P(X_i = 1 | Y = 0) = \theta_{i0} \Rightarrow P(X_i = 0 | Y = 0) = 1 - \theta_{i0}$. Făcând aceste notații putem scrie acum probabilitatea $P(X_i | Y = 1)$ într-un mod compact astfel:

$$P(X_i | Y = 1) = \theta_{i1}^{X_i} (1 - \theta_{i1})^{1-X_i} \quad (14)$$

Formula lui Bayes ne permite să exprimăm $P(Y = 1 | X)$ astfel:

$$P(Y = 1 | X) = \frac{P(Y = 1)P(X | Y = 1)}{P(Y = 1)P(X | Y = 1) + P(Y = 0)P(X | Y = 0)} \quad (15)$$

Scotând factor comun $P(Y = 1)P(X | Y = 1)$ la numitor și făcând simplificările obținem:

$$P(Y = 1 | X) = \frac{1}{1 + \frac{P(Y=0)P(X=0)}{P(Y=1)P(X|Y=1)}} \equiv \frac{1}{1 + \exp(\ln(\frac{P(Y=0)P(X=0)}{P(Y=1)P(X|Y=1)}))} \quad (16)$$

Vom nota probabilitățile a priori $P(Y = 1)$ și $P(Y = 0)$ cu π respectiv $1 - \pi$. De asemenea, întrucât Bayes Naiv lucrează cu presupunerea că pentru un Y dat, $\forall i$ și $j \neq i$, $X_{\S i}$ sunt X_j independente mai

putem rescrie ecuația de sus astfel:

$$P(Y = 1 | X) = \frac{1}{1 + \exp(\ln(\frac{P(Y=0)}{P(Y=1)} + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}))} \equiv \frac{1}{1 + \exp(\ln(\frac{1-\pi}{\pi} + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}))} \quad (17)$$

În continuare, conform (14), vom scrie $P(X | Y = 1)$ ca $\theta_{i1}^{X_i} (1 - \theta_{i1})^{1-X_i}$ și $P(X_i | Y = 0)$ ca $\theta_{i0}^{X_i} (1 - \theta_{i0})^{1-X_i}$. Ecuația devine:

$$P(Y = 1 | X) = \frac{1}{1 + \exp(\ln \frac{1-\pi}{\pi} + \sum_{i=1}^d \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} + \sum_{i=1}^d X_i (\ln \frac{\theta_{i0}}{\theta_{i1}} - \ln \frac{1-\theta_{i0}}{1-\theta_{i1}}))} \quad (18)$$

Pentru ca aceasta ecuație să aibă forma regresiei logistice, vom alege:

$$w_0 = \ln \frac{1-\pi}{\pi} + \sum_{i=1}^d \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \text{ și } w_i = \ln \frac{\theta_{i0}}{\theta_{i1}} - \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \quad (19)$$

2.3.4 Algoritmul Bayes Naiv pentru variabile continue

Deși formula pentru prezicere a etichetelor folosind algoritmul Bayes Naiv este aceeași indiferent dacă lucrăm cu date cu valori discrete sau continue, ceea ce diferă este modul în care calculăm valoarea $P(X_i | Y)$.

În cazul când avem valori discrete, valoarea $P(X_i | Y)$ se obține imediat folosind estimarea MAP sau MLE. În cazul valorilor continue, cea mai folosită modalitate de a reprezenta distribuția $P(X_i | Y)$ este să se folosească presupunerea că pt fiecare valoare discretă y_k , distribuția lui X_i este Gaussiană și anumită este definită de o medie și deviație standard specifice lui X_i pentru y_k . În continuare, aceste valori le vom reprezenta astfel:

$$\mu_{ik} = E[X_i | Y = y_k] \quad (20)$$

$$\sigma_{ik}^2 = E[(X_i - \mu_{ik})^2 | Y = y_k] \quad (21)$$

De asemenea, trebuie să estimăm și probabilitatea a priori a lui Y :

$$\pi_k = P(Y_k = y) \quad (22)$$

Folosind aceste valori putem scrie $P(X_i | Y)$ astfel:

$$P(X_i | Y = y_k) = \frac{1}{\sqrt{2\pi}\sqrt{\sigma_{ik}^2}} e^{\frac{-(X_i - \mu)^2}{2\sigma^2}} \quad (23)$$

Pentru a estima media si deviatia standard putem folosi MLE si obtinem:

$$\hat{\mu}_{ik} = \frac{1}{\sum_j \sigma(Y_j = y_k)} \sum_j X_i^j \sigma(Y^j = y_k) \quad (24)$$

$$\hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \sigma(Y_j = y_k)} \sum_j (X_i^j - \mu_{ik})^2 \sigma(\hat{Y}^j = y_k) \quad (25)$$

2.3.5 Bayes Naiv Gaussian

Folosind aceeasi presupunere ca mai sus si anume ca pentru fiecare X_i , $P(X_i | Y = y_k)$ are o distributie Gaussiană data de parametrii μ_{ik} si σ_{ik}^2 , vom exprima acum termenul $P(Y | X)$. In plus fata de mai sus, mai facem presupunerea ca Y este o variabila booleana, avand o distributie Bernoulli data de parametrul $\pi = P(Y = 1)$. Vom observa ca se va obtine aceeasi forma pe care o folosim la alt algoritm de clasificare: Regresia Logistica (vezi sectiunea 3).

Formula lui Bayes ne permite sa exprimam $P(Y = 1 | X)$ astfel:

$$P(Y = 1 | X) = \frac{P(Y = 1)P(X | Y = 1)}{P(Y = 1)P(X | Y = 1) + P(Y = 0)P(X | Y = 0)} \quad (26)$$

Scoatand factor comun $P(Y = 1)P(X | Y = 1)$ la numitor si facand simplificarile obtinem:

$$P(Y = 1 | X) = \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}} \equiv \frac{1}{1 + \exp(\ln(\frac{P(Y=0)P(X=0)}{P(Y=1)P(X|Y=1)}))} \quad (27)$$

Întrucât Bayes Naiv lucrează cu presupunerea că pentru un Y dat, $\forall i$ și $j \neq i$, X_{Si} sunt X_j independente mai putem rescrie ecuația de sus astfel:

$$P(Y = 1 | X) = \frac{1}{1 + \exp(\ln(\frac{P(Y=0)}{P(Y=1)} + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}))} \equiv \frac{1}{1 + \exp(\ln(\frac{1-\pi}{\pi} + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}))} \quad (28)$$

Utilizând (23) și făcând calculele putem scrie:

$$\sum_i \ln \frac{P(X_i = 0)}{P(X_i | Y = 1)} = \sum_i (X_i \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}) \quad (29)$$

(28), (29) \implies

$$P(Y = 1 | X) = \frac{1}{1 + \exp(\ln \frac{1-\pi}{\pi} + \sum_i (X_i \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}))} \quad (30)$$

Notăm $\ln \frac{1-\pi}{\pi} + \sum_i (\frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2})$ cu w_0 , respectiv $\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2}$ cu $w_i \forall i \geq 1$ și obținem:

$$P(Y = 1 | X) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)} \quad (31)$$

2.4 Concluzie capitol

În acest capitol am făcut o prezentare scurtă a algoritmului Bayes Naiv și modul în care acesta lucrează atunci când avem date discrete și continue. De asemenea, am arătat cum formă folosită de Regresia Logisitica poate fi echivalentă cu cea folosită de Algoritmul Bayes Naiv. Un alt concept introdus a fost o variație al algoritmului Bayes Naiv-algoritmul Bayes Naiv Gaussian-care presupune o distribuție Gaussiană asupra datelor de antrenare.

3 Modelul Discriminativ

3.1 Introducere capitol

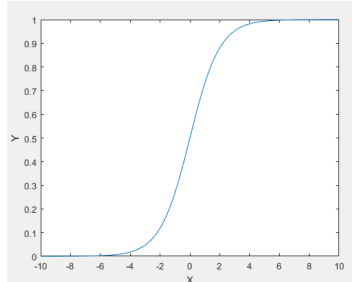
Comparativ cu modelul generativ, un model discriminativ este un model care nu se interesează de modul în care datele au fost generate. Acesta din urmă încearcă să învețe niște granițe de decizie pe care le va folosi mai apoi pentru clasificare unei noi instanțe.

În exemplul de clasificare între motociclete și mașini, un algoritm discriminativ va învăța întâi diferențele dintre o motocicletă și o mașină și folosindu-se de aceste cunoștințe va determina clasa unui nou obiect.

3.2 Regresia Logistica

3.2.1 Introducere Regresie Logistica

Figura 1: Funcția sigmoid



Ca în orice problemă de clasificare, datele de intrare sunt n instanțe de antrenament X_1, \dots, X_n , $X_i = x_1, \dots, x_d \forall i$. Fiecărei instanțe îi corespunde o etichetă Y_i . Regresia încearcă să învețe o funcție care să aproximeze distribuția $P(Y | X)$. Pentru a face acest lucru, acest algoritm presupune aproximată $P(Y | X)$ poate fi ai cu o funcție sigmoidală iar apoi estimează parametrii acestei funcții direct din datele de antrenare.

$$\sigma(x) \stackrel{\text{def}}{=} \frac{1}{1 + e^{-x}} \quad (32)$$

În cazul în care Y este o variabilă booleană, regresia logistică folosește următoarea ecuație pentru a exprima $P(Y = 1 | X)$:

$$P(Y = 1 | X) = \prod_i \sigma(w_i X_i) \quad (33)$$

Nota: din (33) se obține imediat $P(Y = 0 | X)$ folosindu-ne de faptul că suma acestor 2 probabilități

este 1.

În continuare vom arată modul în care regresia logisitică alege valori optime pentru parametrii w_i cum se aleg regulile de actualizare pentru acești parametri.

3.2.2 Estimarea parametrilor

Regresia Logistică alege parametrii w_i urmărind să maximizeze log-verosimilitatea condițională a întregului set de date de antrenament. Verosimilitatea condicională este probabilitatea sa observăm etichetele Y , fiind dat setul de date X . Notăm cu W vectorul $w_0, w_1, w_n \dots$. Ne dorim să alegem W astfel încât:

$$W \leftarrow \operatorname{argmax}_W \prod_l P(Y^l | X^l, W) \quad (34)$$

Observăm că putem scrie probabilitatea $P(Y = y | X = x)$ astfel:

$$P(Y = y | X = x) = \sigma(wx)^y [1 - \sigma(wx)]^{1-y} \quad (35)$$

Folosind această ecuație verosimilitatea condițională devine:

$$\prod_l P(Y = y^{(l)} | X = x^{(l)}) = \prod_l (\sigma(wx^{(l)})^{y^{(l)}} [1 - \sigma(wx^{(l)})]^{1-y^{(l)}}) \quad (36)$$

Aplicând logaritmul și notând log-verosimilitatea condițională cu $l(W)$ obținem:

$$l(W) = \sum_l (y^{(l)} \ln \sigma(wx^{(l)}) + (1 - y^{(l)}) \ln [1 - \sigma(wx^{(l)})]) \quad (37)$$

Întrucât nu există o formulă directă pentru determinarea acelor valori ale lui W care să maximizeze log-verosimilitatea condițională a datelor, vom determina aceste valori folosind o metodă de optimizare. Există mai multe metode de optimizare dar în continuare vom alege o singură metodă: gradientul ascendent. Ideea acestei metode este să pornim cu niște valorilor inițiale pentru parametrii w (valori alese aleator sau în urmă unor cunostine anterioare) și apoi să facem în mod continuu pași mici în direcția gradientului funcției $l(W)$ (vectorul de derivate parțiale). Astfel, vom ajunge într-un final în punctul de maxim global.

3.2.3 Actualizarea valorilor w_i

Pentru a maximiza funcția de log verosimilitate condițională, la fiecare pas optimizăm parametrul w astfel:

$$W_j^{nou} = W_j^{vechi} + \alpha \frac{\delta}{\delta W_j^{vechi}} L(W^{vechi}) \quad (38)$$

unde α este marimea pasului, numita rata de invatare.

În continuare vom calcula valoarea derivatei funcției de log-verosimilitate condițională. Aceasta se va calcula ca o sumă de derivate parțiale. Pentru o instanță (x, y) derivata se calculează astfel:

$$\begin{aligned} \frac{\delta}{\delta W_j} L(w_j) &= \frac{\delta}{\delta W_j} y \ln \sigma(wx) + \frac{\delta}{\delta W_j} (1 - y) \ln [1 - \sigma(wx)] \\ &= \left[\frac{y}{\sigma(wx)} - \frac{1 - y}{1 - \sigma(wx)} \right] \frac{\delta}{\delta w_j} \sigma(wx) \\ &= \left[\frac{y}{\sigma(wx)} - \frac{1 - y}{1 - \sigma(wx)} \right] \sigma(wx) [1 - \sigma(wx)] x_j \\ &= \frac{y - \sigma(wx)}{\sigma(wx) [1 - \sigma(wx)]} \sigma(wx) [1 - \sigma(wx)] x_j \\ &= [y - \sigma(wx)] x_j \end{aligned}$$

Însumând acești termeni pentru fiecare instanță de antrenare obținem derivata log verosimilității condiționale a întregului set de date:

$$\frac{\delta}{\delta w_j} L(w) = \sum_i [y^{(i)} - \sigma(wx^{(i)})] x_j^{(i)} \quad (39)$$

Notă: Mai sus ne-am folosit de următoarea proprietate pentru derivata funcției logistice σ în raport cu argumentul ei:

$$\frac{\delta}{\delta w_j} \sigma(z) = \sigma(z) [1 - \sigma(z)] \forall z \in R \quad (40)$$

3.2.4 Teorema de convergența pentru Regresia Logistica

Pentru a demonstra că algoritmul de optimizare al parametrului W ajunge într-un maxim global vom demonstra că funcția de log-verosimilitate condițională a datelor de antrenare este concavă și deci nu are alt maxim local. Vom scrie întâi matricea gradient H a funcției $L(w)$ formată din pătratele derivatelor de ordin 2 ale parametrului W . Vom arăta apoi că pentru orice vector $z =$

(z_1, z_2, \dots) următoarea ecuație este adevărată: $z^T H z \leq 0$. Acest lucru înseamnă că matricea H este negativ semi-definită și deci funcția $L(w)$ este concavă.

Folosindu-ne de ecuațiile (30), (31) putem exprima termenii matricii H astfel:

$$H_{kl} = \frac{\delta^2 L(w)}{\delta w_k \delta w_l} = - \sum_i \frac{\delta \sigma(w x^{(i)})}{\delta w_l} x_k^{(i)} = - \sum_i \sigma(w x^{(i)}) (1 - \sigma(w x^{(i)})) x_l^{(i)} x_k^{(i)} \quad (41)$$

tiind că $X = x x^T \iff X_{ij} = x_i x_j$ matricea H devine:

$$H = - \sum_i \sigma(w x^{(i)}) (1 - \sigma(w x^{(i)})) x^{(i)} x^{(i)T} \quad (42)$$

Pentru a arata ca H este negativ semi-definita, vom arata ca $z^T H z \leq 0 \forall z$

$$\begin{aligned} z^T H z &= - z^T \left(\sum_{i=1}^m \sigma(x^{(i)}) (1 - \sigma(x^{(i)})) x^{(i)} x^{(i)T} \right) z \\ &= - \sum_{i=1}^m \sigma(x^{(i)}) (1 - \sigma(x^{(i)})) z^T x^{(i)} x^{(i)T} z \\ &= - \sum_{i=1}^m \sigma(x^{(i)}) (1 - \sigma(x^{(i)})) (z^T x^{(i)})^2 \leq 0 \end{aligned}$$

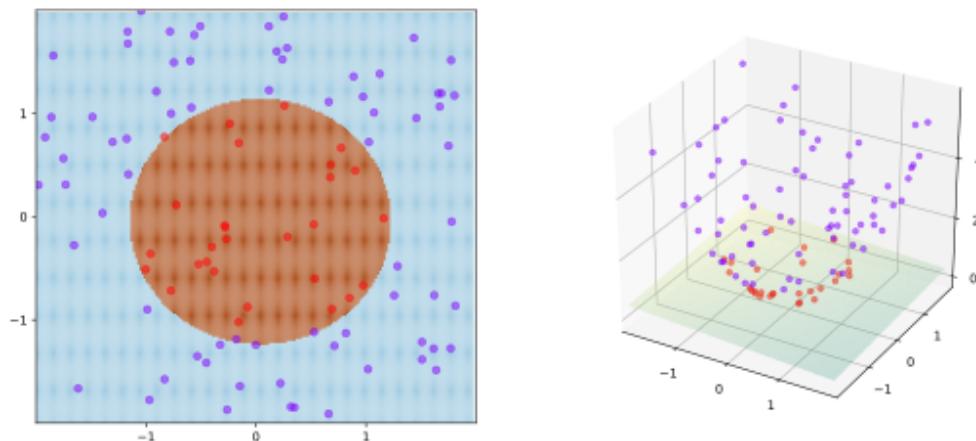
Notă: Ultima inegalitate este adevărată întrucât $0 \leq \sigma(x^{(i)}) \leq 1$ și $(z^T x^{(i)})^2 \geq 0$

3.2.5 Regularizarea regresiei logistice

Atunci când datele de antrenare au multe atribute sau numărul de insante pentru antrenare este redus o problemă care poate apărea folosind regresia logistică este fenomenul de "overfitting". Astfel, algoritmul va avea o acuratețe mare pentru datele de antrenare dar una redusă pentru datele de testare. O modalitate de a evita acest lucru este regularizarea regresiei logistice. Mai exact, vom creă o versiune modificată a funcției de log verosimilitate condițională care va penaliza valorile mari ale parametrului W .

$$W \leftarrow \operatorname{argmax}_W \sum_l \ln(P(Y^l | X^l, W)) - \lambda \|W\|^2 \quad (43)$$

Figura 2: Maparea datelor initiale intr-un spatiu 3d unde se observa ca acestea pot fi separate de un plan. Pentru mapare s-a folosit functia Kernel $K(x, y) = xy + x^2y^2$. (Figura produsa de Shiyu Ji - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=60458994>)



În funcție de modul în care este penalizat paramtrul W , există mai multe tipuri de regularizări. Cele mai întâlnite tipuri sunt regularizarea L1 și regularizarea L2. Un model care folosește regularizarea L1 se numește Regresie Lasso iar un model care folosește regularizarea L2 se va numi Regresie Ridge. Regularizarea L1 penalizează funcția de verosimilitate condițională cu $\lambda \sum_i |w_i|$ iar regularizarea L2 cu $\lambda \sum_i w_i^2$. Deși ambele au ca scop evitarea producerii fenomenului de overfitting există câteva diferențe care ne vor ajuta să alegem tipul de regularizare potrivit. Regularizarea L1 este ineficientă computațional atunci când nu avem multe elemente de 0. De asemenea această va produce multe dintre elementele lui W să fie 0, un lucru avantajos atunci când ne dorim să facem și selectare de trăsături pentru setul nostru de date.

3.2.6 Regresia Logistica Kernelizata

Atunci cand ne confruntam cu un set de date care nu sunt separabile liniar putem incerca sa transpunem datele intr-un spatiu de dimensiune mai mare in care datele vor deveni separabile liniar.

In general, kernelizarea presupune gasirea unei functii de mapare $f : X \rightarrow Z$, unde X este spatiul initial al datelor iar Z este un spatiu de dimensiune mai mare, astfel incat calculele in Z sa se foloseasca doar de produsul scalar al vectorilor si exista o functie K , numita kernel, astfel incat produsul scalar dintre $f(x_i)$ si $f(x_j)$ sa fie $K(x_i, x_j)$.

Notând cu n dimensiunea spațiului de date X și cu m dimensiunea spațiului Z în care ne dorim să proiectăm datele, forma generală a regresiei logistice va deveni:

$$P(Y = 1 | f(X)) = g(w_0 + \sum_{i=1}^m w_i f(X)_i) \quad (44)$$

unde g este funcția sigmoid

Verosimilitatea parametrului W este

$$l(W) = \sum_{l=1}^R Y^l (w_0 + \sum_{j=1}^R w_j (X_j, X)) - \ln(1 + \exp(w_0 + \sum_{j=1}^R w_j (X_j, X))) \quad (45)$$

Pentru a obține regula de actualizare a parametrului W , calculăm derivată ecuației de mai sus:

$$\frac{\delta(l(W))}{\delta w_i} = \sum_{l=1}^R (Y^l - \frac{\exp(w_0 + \sum_{j=1}^R w_j (X_j, X))}{1 + \exp(w_0 + \sum_{j=1}^R w_j (X_j, X))}) K(X_i, X) \quad (46)$$

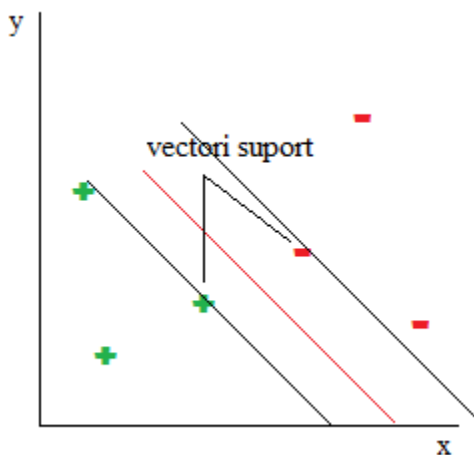
Regula de actualizare:

$$W^{nou} = W^{vechi} + \alpha \frac{\delta(l(W))}{\delta W} \quad (47)$$

3.3 Mașini cu Vector Suport

Un alt algoritm discriminativ care a devenit foarte popular în ultima vreme sunt Mașinile cu vector suport (SVM). Mașinile cu vector suport au ca scop desenarea unui hiperplan care împarte cel mai bine setul de date în clasele în care se dorește a fi făcută clasificarea. .

Figura 3: SVM



Considerăm instanțele $x_1, x_2, \dots, x_m \in R^d$ și etichetele corespunzătoare $y_1, y_2, \dots, y_m \in -1, 1$

Vectorii suport vor fi punctele care sunt cel mai aproape de linia care separă setul de date. Distanța dintre un vector suport și separator poartă numele de margine. În cazul în care datele de antrenare sunt epurabile liniar, scopul algoritmului este să găsească o margine cât mai mare. Găsirea acestei margini este o problemă de optimizare și poartă numele de "Problemă SVM cu margine hard". Odată găsită această margine, folosind figură 3, regulă de decizie este următoarea: noi puncte vor fi clasificate '+' dacă vor

fi la stânga suprafeței de decizie, sau cu '-' în caz contrar.

Matematic, aceasta problema se exprimă astfel:

$$\min_{w, w_0} \frac{1}{2} \|w\|^2 \text{ a.i. } (wx_i + w_0)y_i \geq 1 \forall i = 1..m \quad (48)$$

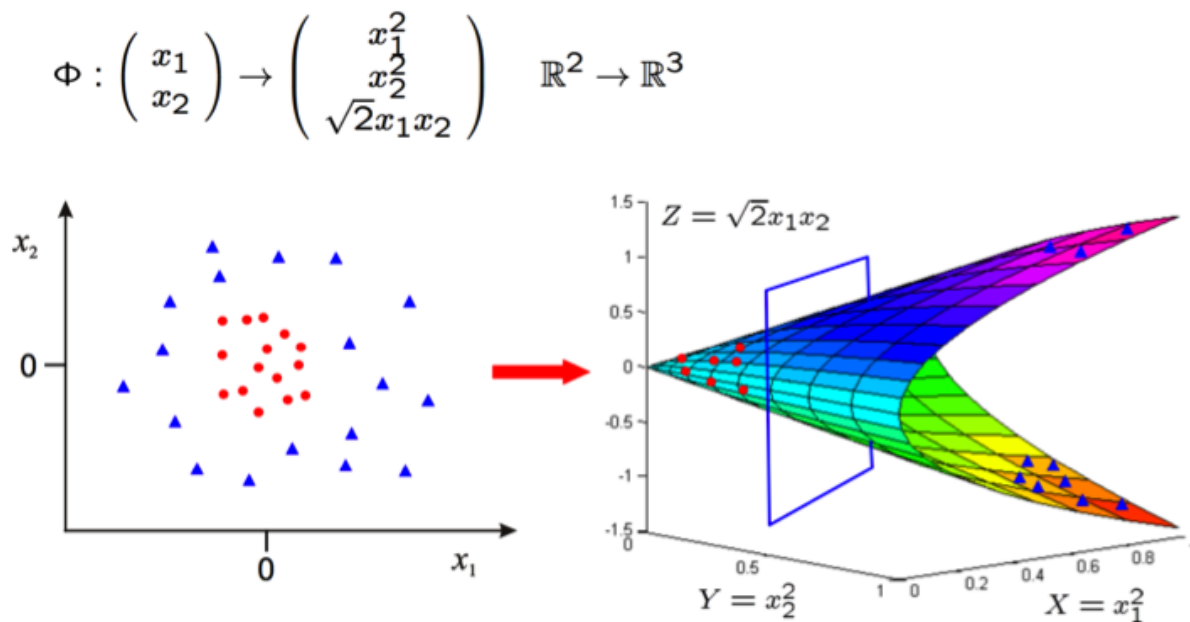
unde $w \in R^d$ și $w_0 \in R$ În urma rezolvării acestei probleme se va obține un model liniar, de forma $y(x) = wx + w_0$, ce va servi ulterior pentru clasificare conform funcției $\text{sign}(y(x))$, echivalenta cu regula de decizie menționată anterior.

În practică însă, datele de antrenare nu sunt întotdeauna perfect separabile liniar. În acest caz folosirea marginii hard nu va mai da rezultate așa că se va folosi "marginea soft". Problemă SVM cu margine soft va conține în plus un parametru de regularizare care va specifică numărul de date clasificate greșit pe care îl permite. Atunci când parametrul de regularizare este mic, clasificatorul va permite că un mic număr din datele de antrenare să fie clasificate greșit. Deși pe datele

de antrenare clasificatorul va avea o performanță destul de bună, acesta nu va generaliza foarte bine, producându-se astfel fenomenul de 'overfitting'. Atunci când termenul de regularizare va fi mai mare, clasificatorul va generaliza mai bine. Totuși, valori prea mari nu sunt indicate întrucât performanță clasificatorului poate scădea destul de mult.

Atunci când datele nu sunt deloc separabile liniar, se va încerca transpunerea datelor într-un spațiu de dimensiune mai mare. Acest lucru este cunoscut sub numele de 'kernelizare'. O funcție kernel este o funcție care cuantifica similaritățile dintre observații.

Figura 4: Datele sunt separabile liniar atunci când le transpunem într-un plan 3D
©towardsdatascience.com



3.4 Concluzie capitol

În acest capitol am descris modul în care un model discriminativ clasifică un set de date, prezentând în special 2 algoritmi: Regresia Loistica și Mașinile cu Vector Suport. Pentru regresia loistica,

care este considerat perechea discriminativa a algoritmului Bayes Naiv, am arătat modul în care se estimează parametrii, precum și 2 metode de îmbunătățire: regularizarea și kernelizarea.

4 Exemple Aplicative

4.1 Aplicare Regresie Logistica

Pentru început vom folosi Regresia Logistică pentru a rezolva o problemă de clasificare binară. Vom aplica algoritmul pe un set de date preluat din baza de date "UCI Machine Learning Repository". Setul de date cuprinde informații despre comunitățile din America precum și o clasificare a acestor comunități dacă sunt sau nu sigure considerând rata infracțiunilor din acea zonă. În total sunt 1994 de observații, fiecare având 128 de atribute.

După ce am rulat algoritmul am obținut o acuratețe de 83.4% . După cum observăm în prima figură, unde am afișat valorile în medie a parametrilor W , raportat la cantitatea în care se găsesc aceștia, obținem valori foarte mari pentru W . Acest lucru nu este de dorit întrucât asta implică o șansa mai mare în apariția fenomenului de overfitting. Pentru a evita acest lucru ne vom folosi de regularizarea L1 și L2 a regresiei logistice.

Figura 1

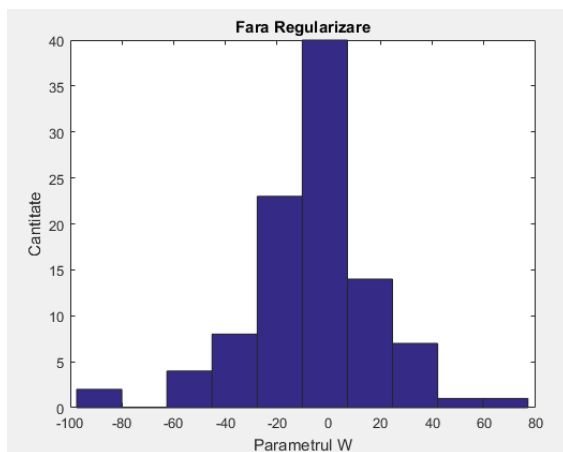
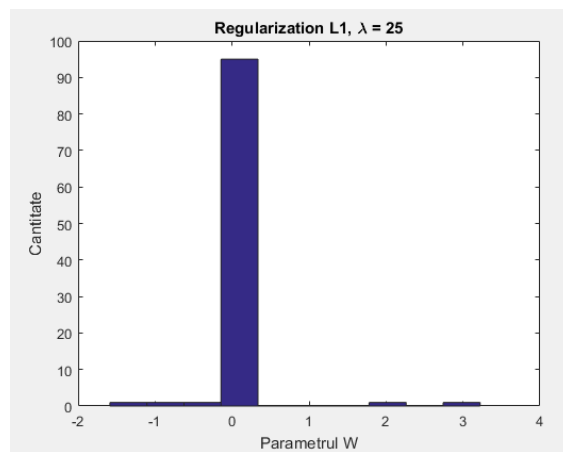


Figura 2 : Regularizare L1

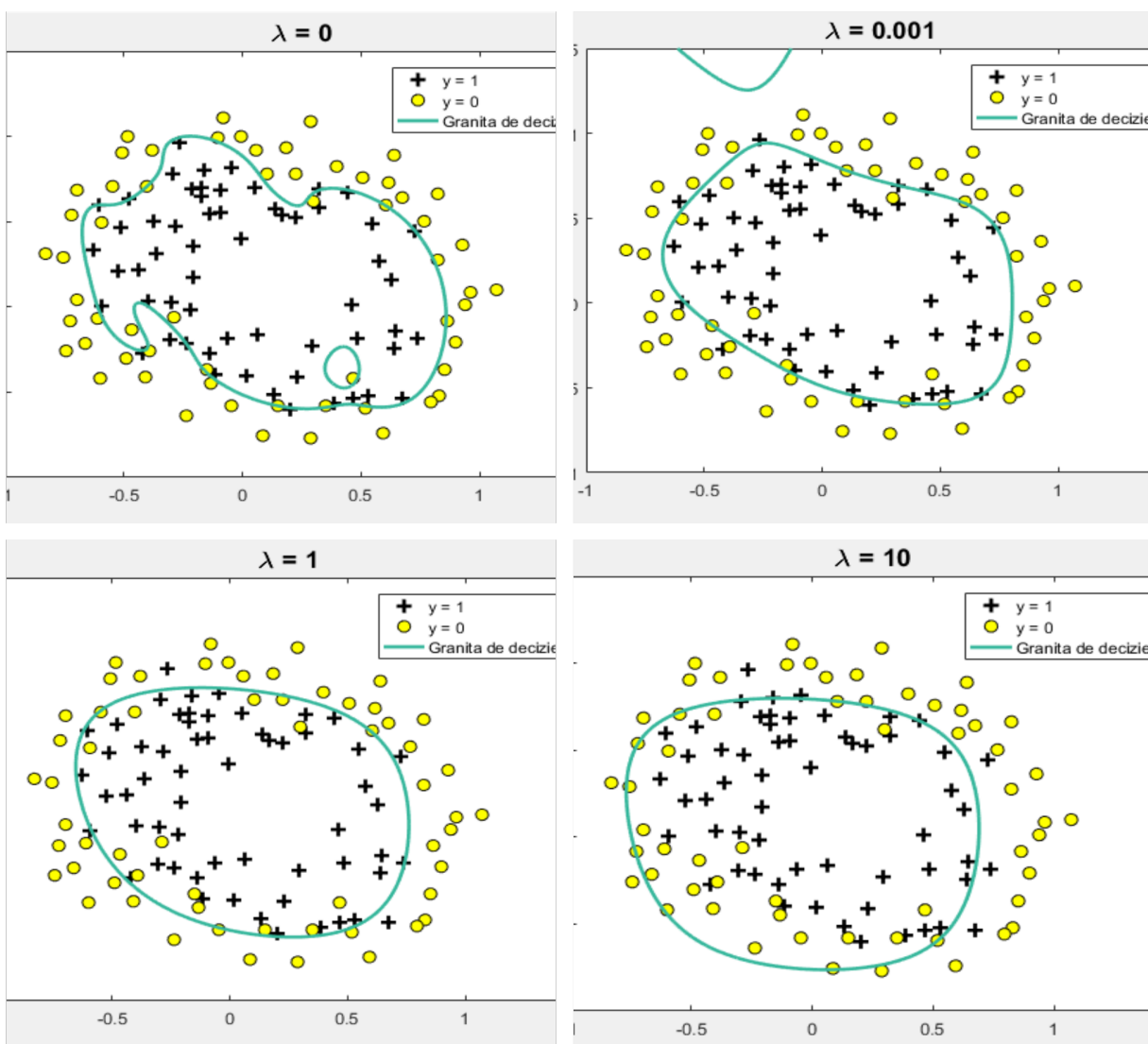


După cum observăm în figura a două, aplicarea regularizării L1 a transformat mulți parametri în 0, ceea ce ne ajută atunci când spațiul de date de antrenare este de o dimensiune foarte mare, întrucât acest lucru ne permite să facem selectare a atributelor. În ceea ce privește acuratețea la cros validare după aplicarea regularizării L1, aceasta nu a crescut foarte mult, procentul exact fiind: 83.5%.

4.1.1 Vizualizarea graniței de decizie

Pentru a vizualiza mai bine modul în care regularizarea regresiei logistice influențează granița de decizie ne vom folosi în continuare de un set de date preluat din cursul de Machine Learning al profesorului Andrew Ng. Datele de antrenare sunt din spațiul R^2 și avem și o clasificare booleană a acestor date. Mai jos am prezentat distribuția inițială a datelor precum și modul în care granița de decizie se modifică în funcție de mărimea parametrului de regularizare λ

Figura 5: Granița de decizie pentru diferite valori λ



După cum observăm în figura de mai sus, atunci când nu avem deloc regularizare algoritmul tinde să se modeleze cât mai mult pe datele de antrenare, ceea ce va determina o eroare mai mare

pe datele de test. În același timp, alegerea unei valori prea mari pentru λ poate diminua cu mult acuratețea clasificatorului atât pe datele de antrenare cât și pe ele de test.

4.1.2 Regresia Logistica MultiClass

Regresia logistică poate fi aplicată și atunci când avem de rezolvat o problemă de clasificare cu mai multe clase. Pentru următoarea aplicație am antrenat regresia logistică pentru a distinge între imagini alb-negru cu cifre. Pentru fiecare cifră din intervalul [0-9] am avut 600 de imagini de antrenare și 500 de testare. Fiecare imagine are dimensiunea de 16*16 px și e reprezentată de un vector de dimensiune 256 în care fiecare element este 0 dacă pixelul respectiv este negru sau 1 dacă pixelul este alb.

După ce am aplicat regresia logistică am obținut o acuratețe de 91.44 %. Aplicând regularizarea L2, am obținut o acuratețe de 91.92% atunci când am ales $\lambda = 10$.

4.1.3 Alte Observatii

Figura 6: Marimea setului de date de antrenare vs Eroare la testare si antrenare ©CMU, T. Mitchel



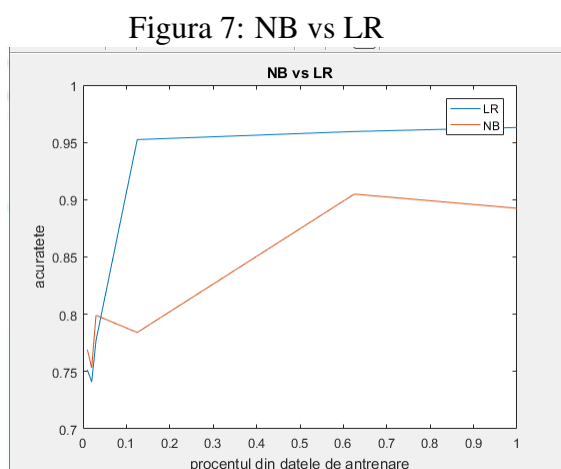
Când setul de date de antrenare este mic, regresia logistică poate de obicei să clasifice perfect datele de antrenare întrucât variant datelor este foarte redusă. Acesta este motivul pentru care eroarea la antrenare este aproape 0. Acest model nu are însă capacitatea de generalizare, întrucât estimările parametrului w sunt bazate pe un subset de date care nu este reprezentativ pentru întregul set de date. Acest lucru este cunoscut sub numele de overfitting. Pe măsură ce avem mai multe date de antrenare, variantă datelor crește și modelul nu mai este capabil de obicei să prezică corect toate datele. În același timp, mai multe date de antrenare oferă modelului nostru o viziune mai largă asupra întregului set de date ceea ce îi permitea aibă o acuratețe mai mare în estimarea parametrului w . Acest lucru duce la o generalizare mai bună, deci o eroare la testare mai mică.

4.2 Aplicare Bayes Naiv vs Regresia Logistica

Pentru a arăta diferența dintre Bayes Naiv și Regresia Logistica am lucrat pe un set de date din baza de date "UCI Machine Learning Repository". Setul de date cuprinde informații despre pacienți diagnosticați cu cancer la sân precum și o clasificare a tumorilor canceroase ca fiind inofensive sau periculoase. În total setul de date are 699 de observații, din care am eliminat 6 pentru care aveam informații lipsă. Fiecare observație are 10 atribute.

Pentru a compara cei doi algoritmi am folosit metoda 10-fold cross validation în care împărțim setul de date în 10 părți egale, antrenăm algoritmul pe 9 subseturi de date și testăm pe un sub-set. Acest lucru îl facem pentru fiecare din cele 10 subseturi. În urma cross validării am obținut următoarele rezultate:

- Bayes Naiv: 0.9103 acuratete
- Regresia Logistica: 0.9853 acuratete



În lucrarea "On discriminative and generative classifiers", profesorul Andrew Ng. spune ca "Pe masura ce marimea setului de date de antrenament crește, initial algoritmul Bayes Naiv va da rezultate mai bune dar in final Regresia Logistica va ajunge la aceeași performanță iar mai apoi probabil va avea chiar rezultate mai bune". Pentru a testa acest lucru am împartit datele în 3 sferturi pentru antrenare și un sfert pentru testare. Apoi, din datele de antrenare am antrenat algoritmul doar cu 1%, 2%, 3%, 12.5%, 62.5% respectiv 100% din setul de antrenament.

Atunci când datele de antrenare sunt într-un număr mai redus, observăm că am obținut o performanță mai mare cu Bayes Naiv: 0.755, comparat cu Regresia Logistica: 0.750. Regresia Logistica a ajuns destul de rapid din urmă însă algoritmul Bayes Naiv, ca în final performanța acestuia să fie mai bună: 0.98 vs 0.91.

4.3 Concluzii Capitol

Dupa cum am aratat mai sus, atunci cand ne confruntam cu o problema de clasificare exista doua moduri de abordare: estimarea directa a parametrilor probabilitatii $P(Y | X)$ - ceea ce numim model discriminativ - sau estimare parametrilor pentru $P(y)$ si $P(X | Y)$, acesta din urma purtand numele de model generativ.

De asemenea, am mai aratat că făcând presupunera că $P(X | Y)$ urmează o distribuție Gaussiană am ajuns la formă probabilității $P(Y | X)$ care este folosită la Regresia Logistica. În continuare, pentru a prezenta câteva avantaje și dezavantaje pentru modelul generativ, respectiv discriminativ vom folosi ca exemple algoritmul generativ Bayes Naiv Gaussian (GNB), respectiv algoritmul discriminativ Regresia Logistică, ambele prezentate mai sus în cadrul acestei lucrări.

Atunci când presupunerile făcute de algoritmul Bayes Naiv Gaussian nu sunt îndeplinite, Regresia Logistica și GNB învață funcții diferite de clasificare. În acest caz eroarea asimptotică a Regresiei Logistice este mai mică decât în cazul algoritmului Bayes Naiv Gaussian. Pe de altă parte, GNB și Regresia Logistică converg la eroarea lor asimptotică la rate diferite. După cum a fost demonstrat într-o lucrare publicată de Ng Jordan (2002), pe măsură ce numărul de instanțe de antrenare n crește, inițial GNB va da rezultate mai bune dar Regresia Logistică va ajunge la un moment dat să aibă acuratețe mai mare.

Concluzionand:

Modelul Generativ	Modelul Discriminativ
Avantaje:	Avantaje:
<ul style="list-style-type: none">• Ne oferă o viziune asupra distribuției datelor• Converge la eroarea asimptotica setul de rapid	<ul style="list-style-type: none">• Nu depinde de presupunerea de independența condițională a datelor• Ușor de modelat
Dezavantaje:	Dezavantaje:
<ul style="list-style-type: none">• Are nevoie de mulți parametri	<ul style="list-style-type: none">• Nu este potrivit atunci când dorim să știm cum au fost generate datele

5 Prezicerea Cancerului la san

5.1 Introducere

În România mortalitatea cauzată de cancerul la sân este de 36%, ceea ce este mult peste media mondială. Conform unei statistici, 90% din cazurile noi de cancer sunt diagnosticate foarte târziu, acestea fiind deja în stadiile ÎI, III sau IV. Deși o parte din vina o poartă pacienții pentru că merg prea târziu la medic, o altă cauza de alarmare este numărul redus de resurse pe care România îl alocă pentru această problemă. Conform unui doctor oncolog din România, dr Dan Jinga: ”În momentul de față suntem doar 360 de oncologi în România. Acești 360 de oncologi împreună cu 140 de hematologi, adică 500 de oameni sunt implicați în administrarea tratamentelor pentru 105.000 de pacienți în fiecare an. Fiecărui medic îi revine în România în medie 450 de proceduri terapeutice pe lună, enorm. Din punctul de vedere al oncologilor, abia avem timp să tratăm pacienții bolnavi.” În aceste condiții, un sistem automat ar fi foarte benefic pentru ajutarea medicilor în diagnosticarea acestei boli. În plus, adunând destule date de antrenare și obținând o acuratețe bună, acest sistem va elimina nevoia ca pacienții să mai treacă prin numeroase teste medicale, care expun pacienții la dureri și radiații.

De-a lungul timpului, au fost propuși mai mulți candidați pentru biomarkerii cancerului la sân. În 1994, a fost desfășurat un studiu ”Machine Learning techniques to diagnose breast cancer from image-processed nuclear features of fine needle aspirates”. În acest studiu au fost analizate 30 de attribute extrase din imaginile scanate. Pe aceste date a fost aplicat un algoritm numit ”Multi-surface Method” care încearcă să plaseze o serie de hiperplanuri în spațiul dimensiunilor de date, încercând să minimizeze numarul de hiperplanuri. Această tehnică a reușit să obțină o acuratețe între 95% și 98%. În 2012, au fost propuși alți 12 biomarkeri(Osteopontin, Haptoglobina, CA15-3, Antigenul Carcinoembrionic, Antigenul Cancer 123, Prolactină, Antigenul Cancer 19-9, *alpha*-fetoproteina, leptina și factorul migrării) însă nu s-a reușit colectarea unor observații semnificative pentru a face o prezicere bună. O cantitate mare de muncă a fost depusă în analiză seturilor de date făcute publice în baza de date ”UCI Machine Learning Repository”: ”Wisconsin Breast Cancer Data Set”, ”Wisconsin Diagnosis Breast Cancer” și ”Wisconsin Prognosis Breast Cancer”. (Bibliografie). Aceste seturi de date au fost colectate de la pacienți în anul 1995. Cele mai multe lucrări

au experimentat pe aceste seturi de date cu diferiți algoritmi de învățare automată, printre care cei mai folosiți au fost: arbori de decizie, rețele neuronale și Mașini cu vector suport.

În continuare, vom aplica algoritmi descriși în această lucrare pe un set de date colectat în 2018 și vom compara modul în care cele două regimuri de abordare, discriminativ respectiv generativ, influențează modul în care algoritmul face o predicție.

5.2 Date Utilizate

Setul de date a fost obținut din baza de date "UCI Machine Learning Repository" și au fost colectate de la un spital din Coimbra. În total sunt 116 de observații, cu 64 de cazuri pacienți bolnavi și 52 sănătoși. Fiecare rând conține 9 atribute diferite și o clasificare: 1 sau 2 reprezentând un pacient sănătos respectiv bolnav. Cele 9 atribute reprezintă date care pot fi colectate din simple analize de sânge și anume: Vârstă, Index de greutate, nivel glucoză, nivel insulină, indice HOMA, nivel leptina, adiponectina, rezistina și nivelul proteinei MCP-1.

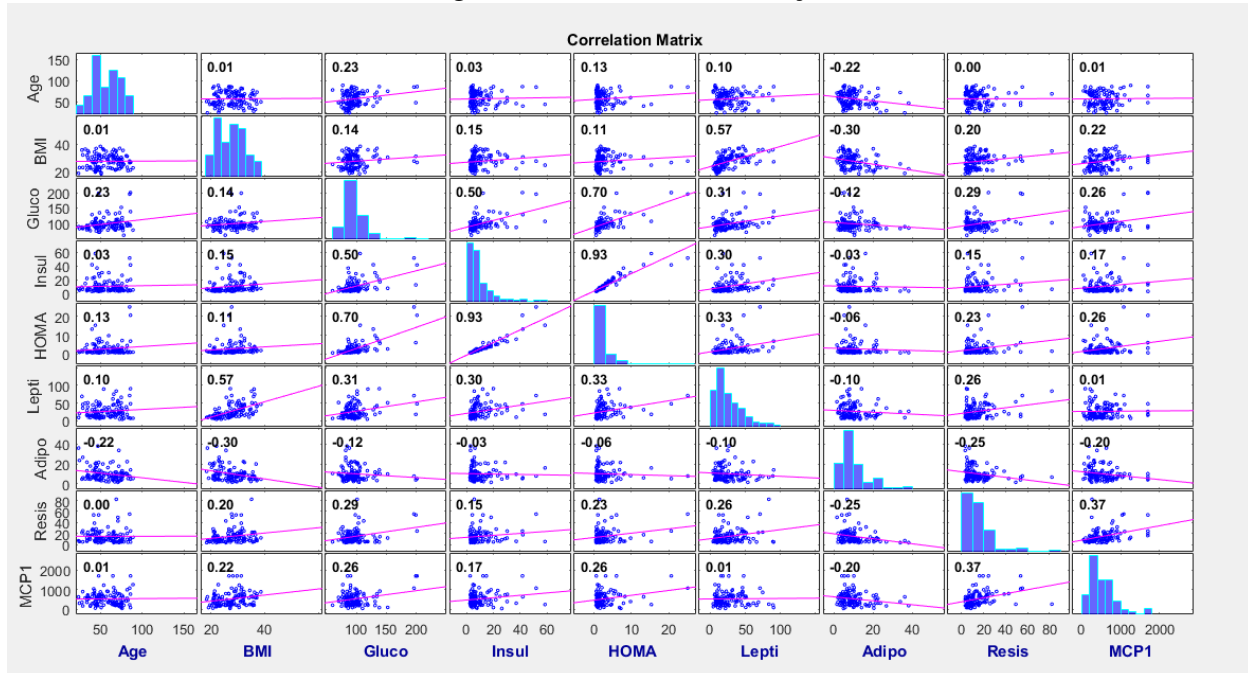
Întrucât setul de date are o dimensiune destul de mică, pentru a testa acuratețea algoritmilor vom aplica cross validarea cu părțile [3, 5, 7, 10]. Astfel, pentru fiecare număr n din mulțimea menționată, vom împărți setul de date în n părți egale, vom testa algoritmul pe una din părți și cu restul vom face antrenarea.

5.2.1 Alegerea Atributelor

În dorință de a obține o acuratețe cât mai mare, primul pas va fi să alegem cele mai semnificative atribute din cele 9. Vom încerca să obținem un număr de 3 sau 4 atribute, ceea ce este mult mai potrivit pentru un set de date așa mic. Pentru a alege atributele, vom aplica mai multe metode iar în final vom alege atributele care sunt clasificate în top 4 trăsături cel mai des.

1) Matricea de Corelație Mai întâi am calculat matricea de corelație dintre variabile, disponibilă în figura de mai jos. Un index de corelație mai mare indică faptul că două atribute sunt destul de asemănătoare și deci alegerea unui singur atribut din cele două este de ajuns. După cum putem

Figura 8: Matricea de Corelație



observa în figura de mai sus, indicele HOMA și nivelul de Insulină ar putea fi grupate împreună, acestea având un nivel de corelație de 0.93.

2) În continuare vom aplica o serie de algoritmi din biblioteca "Feature Selection Library", distribuită open-source în cadrul platformei Matlab.

- Infinite Latent Feature Selection: acest algoritm construiește un graf având ca noduri atributele iar muchiile modelează relația dintre trăsături, mai exact verosimilitatea ca nodurile x și y să fie candidați buni pentru selecție. În urma aplicării acestui algoritm am obținut următoarea ordine a trăsăturilor în funcție de importanța lor: [nivel leptina, nivel glucoză, Vârstă, nivel insulină, indice HOMA, adiponectina, MCP-1, Index de greutate, rezistina]
- Recursive Feature Elimination: acest algoritm începe cu setul inițial de atribute și apoi elimină pe rând câte un atribut care este mai puțin important până se ajunge la numărul dorit de atribute. În urma aplicării acestui algoritm am obținut: [Vârstă, Index de greutate, nivel glucoză, nivel insulină, indice HOMA, nivel leptina, adiponectina, rezistina, MCP-1]

- Scorul Laplacin: acest algoritm se bazează pe observația că în problemele reale de clasificare, observațiile din aceeași clasă sunt în general apropiate una de alta. În urma aplicării acestui algoritm am obținut: [MCP-1, nivel glucoză, rezistina, nivel leptina, Vârstă, Index de greutate, nivel insulină, adiponectina, indice HOMA]
- Scorul Fisher: în acest algoritm, atributele sunt ordanate după un scor calculat pt fiecare atribut în $parte$:

$$scor(s_i) = \frac{\sum n_j (\mu_{ij} - \mu_i)^2}{\sum n_j * \sigma_{ij}^2} \quad (49)$$

unde μ_{ij} și σ_{ij} sunt media și varianța atributului în $clasa$ instanțe în j , n_j este numărul de clasă i în j iar μ_{urma} este media observației glucoză. insulină Vârstă acestui algoritm am obținut: [nivel , indice HOMA, nivel , Index de greutate, rezistina, , MCP-1, nivel leptina, adiponectina]

După ce am rulat acești algoritmi am observat că atributele Vârstă, index de Greutate și index Glucoză apar cel mai frecvent în top 4 atribute.

5.3 Metode și Rezultate

În continuare vom folosi algoritmi decriși în această lucrare pentru a găsi un model care să descrie setul de date. Astfel, vom rula întâi un algoritm generativ: Bayes Naiv, împreună cu variația lui Bayes Naiv Gaussian, iar apoi 2 algoritmi discriminativi: Regresia Logistică (folosind metoda gradientului ascendend și metoda celor mai mici pătrate) și Mașinile cu Vector Suport. Pentru a testa modelul, vom folosi cross-validarea pe 4 niveluri: 3, 5, 7, 10. Astfel, pentru fiecare nivel n , vom împărți setul de date în n părți egale și vom testa pe o singură parte, după ce am antrenat pe restul din setul de date. Fiecare algoritm a fost rulat de 100 de ori, de fiecare dată alegând observații random pentru datele de antrenare respectiv testare. În final, am calculat media preciziilor obținute.

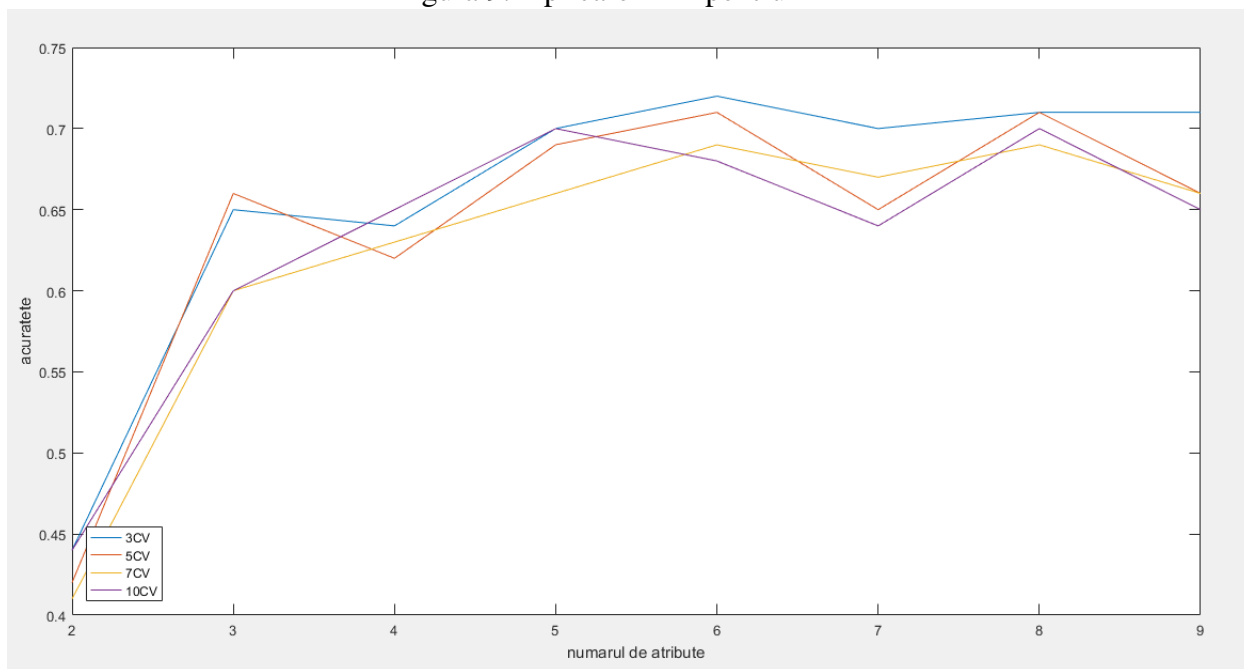
În tabelul de mai jos am pus rezultatele obținute înainte de a aplica un algoritm de alegere a trasterilor

	3-fold CV	5-fold CV	7-fold CV	10-fold CV
NB	0.51	0.57	0.57	0.57
GNB	0.63	0.62	0.63	0.62
LR IRLS	0.71	0.7	0.74	0.73
LR Gr. des.	0.7	0.61	0.55	0.56
LR Kernelizata	0.71	0.66	0.66	0.65
SVM	0.76	0.75	0.7	0.7

În urma aplicării algoritmilor de selecție am rulat din nou algoritmi de mai sus, de data aceasta folosind doar atributele: Vârstă, Index Greutate, nivel glucoză și rezistina.

	3-fold CV	5-fold CV	7-fold CV	10-fold CV
NB	0.66	0.63	0.63	0.63
GNB	0.66	0.7	0.68	0.67
LR IRLS	0.75	0.76	0.75	0.72
LR Gr. des.	0.55	0.56	0.56	0.57
LR Kernelizata	0.65	0.61	0.64	0.64
SVM	0.82	0.78	0.8	0.82

Figura 9: Aplicare RFE pentru LRK



5.4 Concluzii

Observăm că algoritmul Bayes Naiv a dat rezultate mai bune după aplicarea algoritmului de selecție al atributelor. Acest lucru este din cauză că, după cum am observat în matricea de corelație, avem variabile corelate între ele, ceea ce nu este bine pentru algoritmul Bayes Naiv care face presupunerea că datele sunt independente. De asemenea algoritmul Bayes Naiv Gaussian a dat rezultate mai bune decât Bayes Naiv, deci sunt șanse că presupunerea că datele să urmeze o distribuție Gaussiană să fie adevărată.

Ceea ce este interesant este faptul că obținem rezultate mai bune folosind clasificatorii discriminativi, deși setul de date este relativ mic. De asemenea, observăm că obținem rezultate mai bune folosind Regresia Logistică cu metoda iterativă a celor mai mici pătrate (IRLS) decât cu metodata gradientului ascendent. Acest lucru poate fi din cauza că metodă IRLS deși este mai complexă și are nevoie de mai multă putere computațională, este mai rapidă. Astfel, după 100 de iterații regresia logistică cu metooda IRLS a ajuns la o acuratețe mai bună decât regresia logistică folosind gradientul ascendent după 5000 de iterații. Pentru a obține rezultate mai bune și cu regresia logistică folosind gradientul ascendent algoritmul ar trebui rulat de mult mai multe ori.

Deși Regresia Logistică IRLS a avut rezultate similare cu algoritmul Mașini cu Vector Suport înainte de selecția atributelor, după selecție Mașinile cu Vector Suport au crescut semnificativ în acuratețe. Acest lucru ne indică faptul că cele 4 atribute alese de selecție au fost transpuse într-un spațiu dimensional în care acestea sunt relativ separabile liniar.

5.5 Direcții Viitoare

Alegerea unor trăsături semnificative ne-a ajutat să obținem o acuratețe mai mare, însă acest procent tot nu este suficient de mare pentru a prezice cu succes apariția cancerului la sân la un pacient. Întrucât setul de date a conținut doar 116 observații, probabil un set de date mai mare va duce la o performanță mai bună.

6 Bibliografie

1. T. Mitchell, Machine Learning, Chapter 2: Estimating Probabilities: http://www.cs.cmu.edu/~tom/mlbook/Joint_MLE_MAP.pdf
2. T. Mitchell, Machine Learning, Chapter 3: GENERATIVE AND DISCRIMINATIVE CLASSIFIERS: NAIVE BAYES AND LOGISTIC REGRESSION: <http://www.cs.cmu.edu/~tom/mlbook/NBayesLogReg.pdf>
3. Andrew Ng. Michael Jordan: On discriminative vs Generative classifiers: A comparison of logistic regression and Naive Bayes
4. Liviu Ciortuz et al: Exercitii de invatare automata
5. Liviu Ciortuz, curs 6: Bayesian Learning: <https://profs.info.uaic.ro/~ciortuz/SLIDES/ml6.pdf>
6. Tony Jebara, Discriminative, Generative and Imitative Learning: <http://www.cs.columbia.edu/~jebara/papers/jebara4.pdf>
7. Julia Aurelie Lassere, Hybrids of Generative and Discriminative Methods for Machine Learning <http://mi.eng.cam.ac.uk/~jal62/publis/thesis.pdf>
8. Eric Xing, CMU, Lecture 5: <https://www.cs.cmu.edu/~epxing/Class/10701-10s/Lecture/lecture5.pdf>
9. Jing Hao Xue, Aspects of Generative and Discriminative Classifiers: <http://theses.gla.ac.uk/272/1/2008xuephd.pdf>
10. Pattern Recognition and Machine Learning - Christopher M. Bishop
11. Pedro Alonso, A comparison between some discriminative and generative classifiers
12. Andrew Ng, Stanford CS229 Lecture notes, Part IV: <http://cs229.stanford.edu/notes/cs229-notes2.pdf>

13. Using Resistin, glucose, age and BMI to predict the presence of breast cancer <https://bmccancer.biomedcentral.com/articles/10.1186/s12885-017-3877-1>
14. EduMedical: <https://www.edumedical.ro/prevalenta-in-cazurile-de-nouline-cancer-la-san-a-crescut-enorm-in-romania/>
15. Masini cu Vector Suport: <https://towardsdatascience.com/support-vector-newline-machines-a-brief-overview-37e018ae310f>