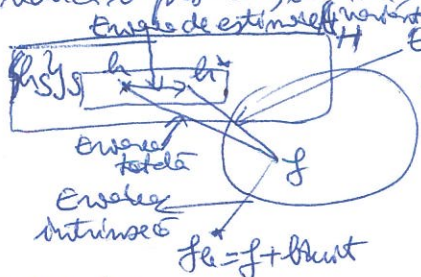


Minimare supravizată

-1-

[CURS]

Compromisul deplasare-dispersie; f (din F = familie teoretică a tuturor funcțiilor posibile care generalizează datele) este o funcție de dependență a valorii introduse x și a etichetei y ; h (din H = spațiul ipoteză) model de funcții, în care câțiva (putem) este o funcție care se dă drept cea mai bună aproximație de f ; h^* = funcție optimă din H (cea mai apropiată de f). Exemplu de minimare supravizată: $S = \{z_1 = (x_1, y_1), \dots, z_m = (x_m, y_m)\}$. h^* = funcție optimă obținută în baza minimării supravizate S , w = pondere; deplasare = diferența dintre F și H , din care H nu conține f . Varianța = deplasare (putem să-l estimăm) în jurul soluției ideale h^* și este adătită cu flexibilitatea spațiului H . Varianța mare = predicțiile vor fi mult de la esențion la actual. Overfitting: modelul minimă h , particularitățile intrinsecă (not-se-ul) ale esenționului de antrenament, nu descrie starea reală generată de f . Generalizare slabă: performanță bună pe datele de antrenament, dar slabă pe date noi. Exemplu/erori intrinsecă = erori de măsurare, transmitere, unare etc. f ipoteză = f + not-se; Compromisul deplasare-dispersie este algebră unde H este diferent de flexibilitate, deci se reduce deplasarea, dar nu atât de bogat încât să explodeze dispersia.



Erora de aproximare (bias) în esență: pentru a reduce deplasarea, H trebuie să fie suficient de mare. Cu cât H e mai mare, cu atât crește deplasarea și erorile totale implicit. În contrast, oare dintre deplasare și dispersie în parte se reduc f mult; Risc empiric ($R_{emp}(h)$) = performanță pe datele de antrenament, $R_{emp}(h) = \frac{1}{|S|} \sum_{i=1}^{|S|} \ell(u_i, h(x_i, w))$. funcție cost (loss) măsoară costul pe care-l plătim în funcție de datele $h(x_i)$ în loc de u_i .

Principiul ERM

găsim h^* care minimizează riscul empiric (eroare de antrenare)

$$h^*_{S,H} = \arg \min_{h \in H} R_{emp}(h); R_{emp}(h) = \frac{1}{|S|} \sum_{i=1}^{|S|} \ell(u_i, h(x_i, w)).$$

afte parametrul w , care minimizează riscul empiric folosind funcție de cost

$$\ell(u_i, h(x_i, w)) = \sum_{i=1}^{|S|} (u_i - h(x_i, w))^2 \quad \text{Funcție de cost generică } \ell(y, a) = \begin{cases} 0, & \text{dacă } a=y, \\ 1, & \text{dacă } a \neq y \end{cases}$$

$$\ell(u_i, h(x_i, w)) = \sum_{i=1}^{|S|} (u_i - h(x_i, w))^2$$

$$\ell(u_i, h(x_i, w)) = \sum_{i=1}^{|S|} |u_i - h(x_i, w)|$$

Risc real = performanță estimată pe orice date noi din aceeași distribuție \rightarrow trebuie redus la maximum, dar pot sărit cel mai aproape de F . Dacă toate erorile au aceeași pondere, atunci riscul real estimat = (suma feloim se calculează doar pe S (setul de antrenament), este t). Risc empiric estimat, nu real. Penalizare = adăugarea costului/găzduirii costului total de reale. \rightarrow

ERM = empirical risk minimization (selectarea modelului care dă cea mai mică eroare pe setul de antrenament). Structural Risk Minimization (SRM)

$$R_{real}(h) \leq R_{emp}(h) + \text{Complexity Penalty}(H, m, \delta) \quad \text{unde } m = \text{numărul de exemple din setul de antrenament}$$

Definiție: Problema modelului de selecție constă în alegerea, în baza unui esențion de date de antrenament, de lungime m , a clasei de funcții optime h din H , a. d. riscul real $R_{real}(h)$ e minim.

Problema estimării riscului real din riscul empiric se reduce la a-i adăuga o penalizare de complexitate care:

1. Coste adătită cu măsura "sau flexibilitate" lui H (relativ constant sau ușor accelerat). În același timp, riscul empiric scade la început mai rapid, apoi din ce în ce mai lent, astfel că suma va forma un grafic în formă de U.
 2. Scade adătită cu numărul de exemple de antrenament m
 3. Separați de nivelul de încredere (delta mic). În practică, pentru experimente statistice adesea se ia $\delta \in \{0,05, 0,01, 0,001\} \Leftrightarrow$ % de încredere = $1 - \delta$
- Metode propuse de rezolvare a calculului riscului real, pornind de la cel empiric (SRM):
- The regularization theory (adăugă un termen de penalizare direct în funcție obiectivă (loss) pe care o optimizăm)
 - Minimum Description Length Principle (MDLP) (oricare model plus date poate fi codificat)

ce un set de date. Modelul bun e acela care imprimă cu datele care cea mai mică lungime totală),
 - Akaike Information Criterion (AIC) - (AIC oferă o aproximație a predicției de înfișare, atunci când modelul tău este folosit pe o aproximație adecvată distribuției de generare; alegi modelul cu cea mai mică valoare AIC: modele cu multe parametri sunt taxate, iar modelele care "pred" bine datele sunt războțite).

- Validare încrucișată (împarti setul de date în k subseturi de dimensiuni (aproxim.) egale. Antrenezi modelul pe toate subseturile în afară de cel selectat și evaluezi eroarea pe acesta. Calculezi și media erorilor obținute pe cele k subseturi \Rightarrow estimare a riscului real; Avantaj: folosește eficient toate datele pentru antrenament și test; Problema este costul computațional mare și, de obicei N (de câte ori trebuie antrenat 1 subset) este între 5-10.

$$\hat{R}_{real}(h) = \frac{1}{N} \sum_{i=1}^N \hat{R}_{real,i}(h)$$

O variantă extremă a validării încrucișate este "leave-one-out" (leave-one-out), și este utilă când setul de date este foarte mic. Păcat, se potrivește un exemplu pe test și se antrenează individual pe fiecare din toate celelalte exemple. Rezultă un H cu varietate mai mică, este mai simplu de implementat și mai rapid, dar estimează ratele de eroare foarte mult decât în metode standard, cu un bias mare.

- Bootstrapping (din setul original de m exemple, generezi B noi seturi de antrenament și, extragi aleatoriu m exemple din total (unele de mai multe ori, unele niciodată). Pentru fiecare set eșantion: antrenezi modelul, evaluezi performanța lui pe exemplele care nu au fost selectate (în medie, aproximativ 36,8% din date) și calculezi media erorilor pe cele B eșantioane noi. \Rightarrow estimare a riscului real; Avantaj: funcționează bine și cu date putine, estimează biasul și varianța estimatorului; Problema este costul computațional mare și, poate fi instabil dacă distribuția e dezechilibrată).

- Media Bayesiana (considerăm fiecare model ca o ipoteză cu o probabilitate controlă, apoi o ajustăm în sus sau în jos în funcție de cât de bine se potrivește pe datele de antrenare. Această probabilitate este distribuită de obicei gaussiană, cu o prioritate la predicțiile de ansamblu).

- Metode bag (generăm B seturi de antrenament, cu m elemente prin extragere repetată, antrenăm B modele separat pe fiecare set și facem media aritmetică a rezultatelor modelelor).

- Metode boosting (folosim modele de antrenare pentru fiecare exemplu la o dată, de obicei 1 sau 2; folosim T-a o fel de proces ("weak learners"), care creează modelele precedente; la t -le antrenăm pe fiecare; calculăm erorile = suma probabilității a erorilor precedente; calculăm ponderile modelelor $= \frac{1}{2} * \ln\left(\frac{1-e}{e}\right)$ și combinăm modelul cu ponderile ca mai mare; actualizăm ponderile exemplilor astfel încât setele ca egale și normalizăm distribuția pentru a menține suma ponderilor $= 1$; Rezultat: bias mic, varianță controlată, dar mai complex).

Categorii de metode:

- selecția modelelor: definiți o ierarhie de subspații cu complexitate (crescând) $H_0 \subseteq H_1 \subseteq H_2 \subseteq \dots \subseteq H_d \subseteq \dots$, apoi, având un eșantion de învățare, găsim subspațiul optim din care ne alegem o ipoteză finală; metode: penalizare a complexității (SRM, MDLP), validare prin învățare multiple/re-esantionare (validare încrucișată, bootstrapping). Spațiul funcțiilor polinomiale (curbe) de grad i :

$$H_i = \{h(x, w) = w_0 + w_1 x + w_2 x^2 + \dots + w_i x^i\};$$

$$H_0 \subseteq H_1 \subseteq H_2 \subseteq \dots$$

funcție constantă dreptă parabolă

- tehnici de regularizare: este folosit similar un criteriu de penalizare asociat fiecărui funcție din H care măsoară complexitatea structurii de parametru, fie este de "regularizare" sau derivabilă / minimal local al unei funcții e : derivata de

gradul 1 = 0 si, ceo de gradul 2 > 0 - efortul maxim motivator ptr la optim; punctul de inflexiune e cel care face derivata de gi. 2 = 0, dar in realitate, are sensuri diferite).
- mediu: functia de predictie se formeaza din mai multe functii din H, care primesc ponderi diferite. Astfel, se "netezeste" variatia excesiva a modelelor (metode bayesiana, metode bag) sau costul capacitatii de reprezentare a unui spatiu H neconvex (metodele boosting).

Metode de regularizare

$$l(w; h(x_i; w)) = \sum_{i=1}^N (y_i - h(x_i; w))^2 + \lambda ||w||^2$$

$||w||^2 = w_0^2 + w_1^2 + \dots + w_n^2$ penalizare

λ - controlul importanței termenului de regularizare/penalizare

Alegerea modelului: împartim datele initiale in 2 multimi: multimea de antrenare si multimea de validare si, alegem si λ pe baza erorii pe multimea de validare.

Matricea de confuzie (M_{ij}) - a unei reguli de clasificare este o matrice de dimensiune 2×2 al carui element generic indica numarul de exemple din setul de date T , care apartin clasei i si, care au fost clasificate ca fiind din clasa j . Pn clasificare binara:

	+	-
+	true positives	false positives
-	false negatives	true negatives

TPR (true positive rate / sensibilitate / recall) = TP/P
 FPR (false positive rate / fell-out) = FP/N
 $P = TP + FN$ (toate pozitivele corecte)
 $N = FP + TN$ (-u - negativele corecte)

Deci avem nevoie de rezultate binare si avem mai multe functii h (calculam suma ponderata a lor) sau o functie, ce returneaza valori intr-un interval, folosim 1 prag pentru a genera un rezultat binar. Pntr-o sarcina de detectie binara (obiect prezent versus obiect absent), o formula uzuala de predictie e: $\hat{y}(x) = 1 / (f(x) > \bar{c})$, unde:

- $\hat{y}(x)$ e regula de decizie, care = 1, daca obiectul exista si, 0 daca nu exista;
- $f(x)$ este o functie indicata, care = 1, daca conditia e adevarata si, 0 daca nu e;
- $f(x)$ este scutul de incalzire - poate fi si probabilitatea estimata $P(y=1|x)$, dar, in practica, $f(x)$ vine dintr-un model de genul:

- Reprezentare logistica: $f(x) = \sigma(W^T x + b)$, unde σ = functie sigmoid;
 - SVM: $f(x) = W^T x + b$ (o valoare reala, care se interpreteaza, ca "distanta" fata de hiperplan);
 - Retea neuronală: ultimul strat cu activare softmax sau sigmoid da un scor $f(x) \in [0, 1]$.
- \bar{c} e pragul de sensibilitate a detectiei. In practica, multe modele, elibereaza $f(x)$ ce probabilitate si, iar $\bar{c} = 0,5$.

Variatia valorii pragului / estimarii corecte = curba ROC (caracteristica de operare a receptorului)
 $ROC = (FPR(\bar{c}), TPR(\bar{c}))$, unde \bar{c} = pragul.

And, pt a se gasi un eveniment rar (precum identificarea unui element intr-o imagine), fiindca rata de detectie se intrumpla and pragul \bar{c} e foarte mare (adica sunt putine $f(x) > \bar{c}$, majoritatea lucrurilor sunt false), atunci TPR devine f mic, asa ca se prefera sa se foloseasca pe axa OX: FP in loc de FPR ($FP =$ "cate albne false da detectioul").

Def. Precizie $P = TP / (TP + FP) = p(\hat{y}=1 | \hat{y}=1)$; cat din, ce detectam pozitiv e corect

Def. Recall $R = TP / (TP + FN) = p(\hat{y}=1 | y=1)$; cat % dintre pozitive am detectat

Precizie-recall este un grafic de puncte cu coordonate Recall si Precizie, cu valori calculate pentru fiecare prag. Punctele trebuie sa tinda spre dreapta sus (1 si 1).

Scoala $F1 =$ medie armonică dintre $P(\text{precizie})$ și $R(\text{recall})$ pentru un anumit $\text{flag} =$
 $= 2 \left(\frac{1}{P} + \frac{1}{R} \right)^{-1}$. Acesta este util când este un număr f mic de regulatoare pozitive, iar precizie
 e mai mică decât recall e f mic. Medie aritmetică or f mare, dar medie armonică
 or f mai mică și or evidue mai corect eficiența strategiei. Estimare $R_{\text{real}}(h)$

realizare în funcție de intervalul de încredere, doar numărul de exemple de test este mare
 (> 100):
 $R_{\text{real}}(h) = \left[\frac{\text{terr}}{t} + \zeta(x) \sqrt{\frac{\text{terr}}{t} \left(1 - \frac{\text{terr}}{t} \right)} \right]$. Funcția $\zeta(x)$ (zeta(x)) are ca notă
 valoarea $x =$ numărul de erori
 create din toate testele.

$\rightarrow SE = \frac{\bar{\sigma}}{\sqrt{N}} = \sqrt{\frac{\bar{\sigma}^2}{N}}$, unde $\bar{\sigma}^2$ este estimarea varianței și
 $\bar{\sigma}^2 = \frac{1}{N} \sum_{t=1}^N (L_t - \bar{L})^2$; $L_t = L(y_t, f_m(x_t))$; $\bar{L} = \frac{1}{N} \sum_{t=1}^N L_t$

σ măs. variabilitate intrinsecă a L_t across samples. σ măsura incertitudinii lui L .
 Regula Error Standard: un mod sigur de a alege un model după calcularea riscului și a
 erorii standard, este alegerea celui care corezunde celui mai simplu model al cărui risc
 nu este mai mare decât riscul celui mai bun model + eror standard.

Ptr erorare supraliniară, se folosesc 3 seturi de date mutual exclusive: de antrenament,
 de testare și de validare. Metoda permite și identificarea convergenței dintre diverși algoritmi.
 Compararea a 2 ipoteze obținute cu același algoritm: dacă avem 2 funcții din H care sunt
 obținute de același algoritm pentru seturi de date T_1 și T_2 cu terr_1 și terr_2 numărul de
 teste greșite din totalul t_1 și t_2 , atunci riscul unei dintre funcții față de celelalte e:

$|\hat{\sigma}_R(h_1, h_2)| = \frac{\text{terr}_1}{t_1} - \frac{\text{terr}_2}{t_2}$
 Intervalul de încredere a valorii este dat de formula: $\left[\hat{\sigma}_R \pm \zeta(x) \sqrt{\frac{\text{terr}_1}{t_1} \left(1 - \frac{\text{terr}_1}{t_1} \right) + \frac{\text{terr}_2}{t_2} \left(1 - \frac{\text{terr}_2}{t_2} \right)} \right]$

Compararea a 2 algoritmi (A_1 și A_2), aplicată pe seturi de testare diferite.
 Folosim tehnici ce utilizează selecție aleatorie a seturilor de antrenament și testare, precum validarea
 încrucișată. Exemplu de algoritm:
 împartim setul de date D în k părți egale: T_1, T_2, \dots, T_k
 pentru fiecare k antrenăm pe un subset S_i și testăm pe T_i fiecare algoritm obținând ratele de
 erori (terr_i/t_i) ptr fiecare $R_{p1}, R_{p2}, \dots, R_{pk}$. $\bar{\sigma}_i = R_{p1} - R_{p2}$
 calculăm valoarea medie a celor k valori $\bar{\sigma}$ (sum $\bar{\sigma}_i/k$)
 dacă fiecare parte k are cel puțin 30 de elem. în ea, putem calcula diferența dintre
 intervalele de încredere:

$\left[\bar{\sigma} \pm \zeta(x, k) \sqrt{\frac{1}{k(k-1)} \sum_{i=1}^k (\bar{\sigma}_i - \bar{\sigma})^2} \right]$

Compararea a 2 algoritmi pe același set de testare (testul peche al lui McNemar sau
 Gielick)

Notăm: m_{00} = numărul de exemple clasificate corect de h_1 și h_2
 m_{01} = ————
 m_{10} = ————
 m_{11} = ————
 $2 \geq \frac{|m_{01} - m_{10}|}{\sqrt{m_{01} + m_{10}}}$

Modelul suprimat al învotării din date: - Setul: medonul (setul de date x), profesorul
 (etichetele y asociate x -urilor) și LM.
 - Schema (LM caută în H cea mai bună funcție $h(x)$, care returnează cel mai des
 valoarele y)

- principalul inductiv (induce ce este la un nivel altele obiectivel ideal, ce influențează ^[Curs] ~~obiectivele~~ ^{obiectivele} ~~metodele~~

de învățare)

Inductivitatea este procesul de inducere a ideilor și conceptelor din experiența concretă.

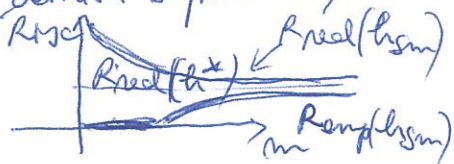
Teoria diferenței inductive. Principul inductiv indică ce prezumție să fie făcută pe a-mbini onto uscul real în funcție de setul de antrenament. Empirical Risk Minimization - (nu este un principiu inductiv - formule de mai sus, se vede uscul real). Alegem

ERM (pp. doc. scade riscuri impozabile) la gestion sau Alegem o singura optiune folosind
folosind principiul de decizie (media) la gestion sau Alegem o singura optiune folosind
principiul compuneri informatiei: eliminam datele redundante ptr a evita o dublare a regulilor
de aplicare a legii. Pe ce aceste sunt valabile si ptr restul datelor.

Condiții de validitate pti principalul CRM: Diferența dintre riscul real pti ipoteze optimă & real (h^*) versus riscul empiric pti ipoteze găsita $R_{emp}(h^*)$ - din setul de antrenament. Considerăm ε o valoare maximă acceptată a diferenței dintre cele 2 riscuri. Diferența depinde de representativitatea datelor de antrenament față de totalul datelor; se va conta funcția care este cea mai bună pti cele mai probabile situații, motiv pti

Cal se plăsește setului de antrenament alse alături
 Constatate prin calculul ERM: Principiul ERM este consistent deo riscul Real
 recunoscând $R_{red}(t, s)$ și, riscul empiric $R_{emp}(t, s)$ converge către acesă, limită, $R_{red}(t, s)$
 Când dimensiunea n a setului de antrenare tinde la ∞ . Riscul empiric este 0 la început

detrital superimposition.



Generalizarea legii numerelor mari (Vapnik & Chervonenkis, 1971, 1989-VC) este absolut necesară și suficientă pentru convergența uniformă nu sunt îndeplinite - adică două entități VC raportate la numărul de observări nu tind către 0 - atunci (\exists) un subspace X^* al spațiului R^n cu măsura de probabilitate egală cu c , cîmădă cumva greșe orice esență de dimensiune arbitrară k , $x_1^*, \dots, x_k^* \in X^*$, poate fi separat în toate cele 2^k moduri posibile de către funcțiile din mulțimea admisibilă de funcții indicator $f(x, q)$, $q \in A_i$.

Aditq:

Fără convergență uniformă ($R_{\text{red}} = R_{\text{test}} \rightarrow 0$, când $m \rightarrow \infty$), nici un algoritm determinist care alege o singură funcție nu poate generaliza în mod garantat. Excepție: dacă introduci informație a priori despre cele funcții sunt "plausibile" (de exemplu un prior bayesian) și poți face mediere / ponderare bayesiană, poți obține, totuși, generalizare. Consecința teoremei: fără nici o informație suplimentară (a priori), ✓ alegere unică de funcție nu "supraviețuiește" doar pe datele de antrenament și / nu va generaliza la date noi: poți găsi întotdeauna un exemplu advers în care acel model se poate vedea perfect (eror 0 pe antrenament), dar e complet eronat pe alte exemple. Entropia VC cuantifică complexitatea clasei de ipoteze; pînă a obț. convergență uniformă, entropia VC la nr. de exemple trebuie să tindă către 0.

entru pe VC la nr. de exemplu trebuie să lăsați ca e.

Protocolul înnotării: stăbilit interacțiunea dintre LM și mediu, incluzând timpul de date de înnotare, timpul de răspuns posibil și supranota. Înnotare batch - toate datele sunt înnotate la început, înnotare on-line - datele sunt înnotate în secvențe, iar LM trebuie să furnizeze un răspuns după fiecare (secvență de identificare explicativă și predictivă).

Modelul boosting - când LM poate influența probabilitatea distribuției de date din care înnotă. Înnotare prin întrebări de categorisire - LM întrebă despre clasele din care aparțin ce a observat.

Ministere active - LM organizează experimente asupra lumii și observă efectele

Probleme legioni vasculare empiric de cel real

a) enunt: modelul de minimizare a riscului: $R(a) = \int Q(y, f(x, a)) dP(y, x), a \in A$, unde probabilitatea nu este cunoscută, dar perechile (y_i, x_i) sunt date.

b. soluții:

b.1. procedura bazată pe gradient: $\alpha_n = \alpha_{n-1} + \delta_n \nabla Q(y_n, f(x_n, \alpha_{n-1}))$, unde α_n = vectorul de parametri la pasul n , δ_n = rata de învățare, ∇ = gradientul (în ce direcție și cu ce intensitate trebuie modificat parametrul), $Q(y_n, f(x_n, \alpha_{n-1}))$ = o funcție care măsoară eroarea (sau pierderea) între valoarea dată y_n și predicția modelului $f(x_n, \alpha_{n-1})$.

b.2. ERM cu condiție convergenței uniforme

c. Demonstrații:

a.1. teorema aproximării stochastice - pentru procedurile bazate pe gradient
a.2. teorema convergenței uniforme pentru principiul ERM.