



La dimension de Vapnik-Chervonenkis

Gilles Richard

Retour sur PAC !

Si C fini et il existe $A...$ alors

- C PAC apprenable... et meme plus car...
- borne sur le nombre d'exemples en $|C|$

Si C infini...

- exemples: $C = \{\text{les rectangles}\}$;-)
- pas de borne generale sur le nb d'exemples

Pb: trouver une autre borne;-)

Idee:

- ds le cas des rectangles, 4nbres suffisent!
- remplacer $\log(|C|)$ par similaire

LA VC-dimension

- C une classe de parties de X (individus)
- $A \subseteq X$, A fini (nos exemples!)

Def1: $\text{tr}(C, A) = \text{tr}_C(A) = \{c \cap A \mid c \in C\} \subseteq 2^A$

Def2: C pulverise A ssi $\text{tr}(C, A) = 2^A$

Def3: $\text{VC-dim}(C) = \text{argmax}\{|A| \mid C \text{ pulverise } A\}$

Ex.:

- Intervalles sur la droite reelle (2)
- Demi-plans dans le plan (3)
- Rectangles //axes: 4 (5: on choisit une classe differente pour le pt qui se trouve a l'int. du rect. defini par les 4 autres)
- Polygones a d cotes: $2d + 1$ (TBD)

Quelques petites propriétés

- $c_a = \{x \text{ in } \mathbb{R} \mid \text{sign}(\sin(ax)) \geq 0\}$:
concept de \mathbb{R}
- $C = \{c_a \mid a \text{ in } \mathbb{R}\}$ (C infini ici)
- $\text{VC-dim}(C) = \infty$
- $|\text{tr}(C, A)| \leq 2^{|m|}$ ou $|A|=m$
- $\text{VC-dim}(C)$ mesure la « puissance d'expression » de C

Prop. 1: Si C finie alors $\text{VC-dim}(C) \leq \log_2(|C|)$

Preuve: le A max est tq : $2^{|A|} \leq |C|$; -)

Un beau resultat!

Prop. 2: VC-dim(C) infinie \rightarrow NON PAC apprenable

Preuve: (1989)

- Idee de la preuve: par l'absurde et on suppose C PAC.
- m nbre d'ex. necessaires pour PAC apprendre avec $\varepsilon=0.1$ et $\delta=0.1$
- On sait que moyenne erreur ≤ 0.19 (vu avant) et
- On sait qu'il existe S pulverisable de cardinal 2^m car VC-dim infinie ;-): on considere ce S comme notre TS
- On considere P uniforme sur TS et 0 ailleurs
- On construit ensuite un concept c tq $P(c(x_i) = 0) = \frac{1}{2}$
- Alors l'erreur moyenne est $\frac{1}{4} > 0.19$: contradiction **CQFD**

Et la reciproque?? Si-dim finie ????

La fonction de croissance Π_C

Notons $\Pi_C(m) = \max\{ |\text{tr}(C,A)| \mid |A| = m \}$

1. Mesure la “complexite/puissance” de C
2. $\Pi_C(m) = 2^{|m|}$ ssi C pulverise un ensemble de card m
3. $\text{VC-dim}(C) = \max \{d \mid \Pi_C(d) = 2^{|d|}\}$
4. $\text{VC-dim}(C)$ infinie ssi $\forall m, \Pi_C(m) = 2^{|m|}$
5. Croissance exponentielle attendue en fonction de m! et pourtant...

Prop. 3: Si $\text{VC-dim}(C)=d$ alors $\Pi_C(m) \leq \Phi_d(m)$

Preuve: (lemme de Sauer – discontinuite ds la fonction de croissance)

- Notons $\Phi_d(m)$ le nbre de parties de moins de d elts ds un ens. a m elements.
- $\Phi_d(m) = \Phi_d(m-1) + \Phi_{d-1}(m-1)$ (explain !)
- D’ou $\Pi_C(m) \leq \Phi_d(m)$ par double induction sur m et d (c’est lourd !)

Prop. 4: $m \leq d$ alors $\Phi_d(m) = 2^{|m|}$ sinon $O(m^d)$

Preuve:

- ($m \leq d$: oui) Si $d < m$, on multiplie par $(d/m)^d \Phi_d(m)$ et on fait rentrer dans la somme car $d/m \leq 1$: en allant jusqu’a m au lieu de d = dlpt en serie de $(1+x)^m$ ou $x = d/m$ qui est inferieur a e^x
- On elimine $(d/m)^d$ de chaque cote ... reste $(em/d)^d$ cqfd

A quoi ca sert?

Notion de ε -reseau S pour $c \in C$ (c target concept):

- $\Delta_\varepsilon(c)$: ensemble des erreurs $r=c\Delta c'$ tq $P(r) \geq \varepsilon$
- r de $\Delta_\varepsilon(c)$, il y a un elt de S dedans
 - un petit exemple $C =$ les intervalles de $[0,1]$
 - $S = k\varepsilon$ avec $1 \leq k \leq \lceil 1/\varepsilon \rceil$ (ie le plus petit entier $\geq 1/\varepsilon$)

Soit TS ε -reseau , soit A algo. polynomial capable de retourner h consistant avec c sur TS

alors $\text{erreur}(h) \leq \varepsilon$ (i.e. $c\Delta h$ n'est pas dans $\Delta_\varepsilon(c)$)

Notons que pour TS de taille m

$p(TS \cap r = \emptyset) \leq (1 - \varepsilon)^m$ donc $p(\text{il existe } r...) \leq |\Delta_\varepsilon(c)| (1 - \varepsilon)^m$

Petite idee de la preuve...

- Event A = "TS n'est pas un ε -reseau "
- But: borner mieux la proba de A : $P(\mu(e) \geq \varepsilon) < P(A)$
- Chercher autre event B tq $\Pr(A) \leq k \Pr(B)$
- Intuitivement augmenter $m = |TS|$;-)))
- On ajoute TS' de taille m a TS
- B_r = "TS n'est pas un ε -reseau \wedge TS' touche r au moins $\varepsilon m/2$ fois (ou r fixe non touche par TS)
- $B = \bigcup B_r$ pour les r non touches par TS
- Of course $P(B) \leq P(A)$ car $B \rightarrow A$!
- Mais le plus beau est: $P(A) \leq 2P(B)$ pour $m = O(1/\varepsilon)$

On y est ... presque

- La notion de trace reapparaît ;-)))
- $P(B) \leq \frac{|\text{tr}_{\Delta\epsilon(c)}(TS \cup TS')|}{P(B_r)}$
 $\leq \frac{|\text{tr}_c(TS \cup TS')|}{2^{-\epsilon m/2}} \leq \Phi_d(2m)$
 $2^{-\epsilon m/2}$
 $\leq (2em/d)^d 2^{-\epsilon m/2}$ (via le lemme de Sauer)
- **Donc** $P(A) = P(\mu(e) \geq \epsilon) \leq 2 (2em/d)^d 2^{-\epsilon m/2}$
- $P(\mu(e) \geq \epsilon) \leq \delta \dots$ la meme chanson!

Et le resultat final

- **Prop.3 : C tq $VC\text{-dim}(C)=d$ finie**

Si $m \geq c_0(1/\varepsilon \log 1/\delta + d/\varepsilon \log(1/\varepsilon))$ alors tt algo A rendant h consistant avec c est PAC.

- **Prop.4: au moins $\Omega(VC\text{-dim}(C)/\varepsilon)$ exemples pour PAC apprendre la classe C avec erreur $\leq \varepsilon$**

Preuve: on construit un TS de taille d pulverise par C, etc...

Conclusion: borne inferieure sur le nbre d'exemples (difficulte statistique!)

V. Vapnik & A. Chervonenkis?

- Education: Russie (Ouzbekistan)
- Now: Prof. London (Royal Holloway)
- Statistical learning theory – SVM
- Chervonenkis... bureau d'en face;-)



Conclusion

- VC-dim: tres importante
- Si finie, strategie de consistance OK
- Relation risque empirique/risque vrai
- Principe de minimisation du risque empirique (ou structurel) (ERM principe)

MAIS

- SVM VC-dim tres large ;-)
- Neural networks idem...
- ILP... infinie ;-)))
autre concept necessaire... probably ;-)

Induction versus Transduction

- **Cadre general (Faire un schema)**
 - Observations $(x_1, y_1), \dots, (x_m, y_m)$
 - Nouveau cas a predire $(x, ?)$
- **Induction: 2 etapes (offline)**
 - particulier \rightarrow general \rightarrow particulier
 - ex.: SVM, neural networks, Bayesian network
- **Transduction: 1 etape (online)**
 - particulier \rightarrow particulier
 - ex.: k-nn, analogical transduction

Limite de PAC...

- **Plus d'hypothese $h...$ donc**
 - Erreur ϵ : encore un sens
 - Confiance δ : pas si clair
 - PAC ... pas si interessant
- **Autre modele?**
 - Gamerman et Vogt (2000)
 - On garde ϵ - on jette δ
- **Cas de la classification**
 - Notion de conformite/non conformite
 - Notion de predicteur conforme