CSCI B609: "Foundations of Data Science"

Lecture 11/12: VC-Dimension and VC-Theorem

Slides at http://grigory.us/data-science-class.html

Grigory Yaroslavtsev <u>http://grigory.us</u>

Intro to ML

- Classification problem
 - Instance space $: \{0,1\}$ or \mathbb{R} (feature vectors)
 - Classification: come up with a mapping $\rightarrow \{0,1\}$
- Formalization:
 - Assume there is a probability distribution over
 - * = "target concept" (set * ⊆ of positive instances)
 - Given labeled i.i.d. samples from $\ \$ produce $\ \ \subseteq$
 - Goal: have agree with * over distribution
 - Minimize: () = $\Pr[\Delta^*]$

Intro to ML

• Training error

 $- = labeled sampled (pairs (,), ∈ , ∈ {0,1})$ - Training error: () = $\frac{| ∩ (Δ *)|}{| |}$

- "Overfitting": low training error, high true error
- Hypothesis classes:
 - H: collection of subsets of called hypotheses
 - If $= \mathbb{R}$ could be all intervals $\{[,], \leq \}$
 - If $= \mathbb{R}$ could be linear separators: $\{\{ \in \mathbb{R} \mid \cdot \geq 0\} \mid \in \mathbb{R}, 0 \in \mathbb{R} \}$
- If is large enough (compared to some property of

Overfitting and Uniform Convergence

• PAC learning (agnostic): For , > 0 if $| | \ge 1/2 \frac{2}{\ln |} + \ln 2 /)$ then with probability 1 - : $\forall \in H: | () - () | \le 1$

- Size of the class of hypotheses can be very large
- Can also be infinite, how to give a bound then?
- We will see ways around this today

VC-dimension

- VC-dim() $\leq \ln |$
- Consider database age vs. salary
- Query: fraction of the overall population with ages 35 45 and salary \$(50 70)K
- How big a database can answer with ± error
- 100 ages × 1000 salaries $\Rightarrow 10^{10}$ rectangles
- $1/2^{2}(10 \ln 10 + \ln 2/)$ samples suffice
- What if we don't want to discretize?

VC-dimension

- **Def.** Concept class shatters a set if $\forall \subseteq$ there is $h \in$ labeling positive and $A \setminus$ negative
- **Def. VC-dim**() = size of the largest shattered set
- Example: axis-parallel rectangles on the plane
 - 4-point diamond is shattered
 - No 5-point set can be shattered
 - VC-dim(axis-parallel rectangles) = 4
- **Def.** $[] = \{h \cap : h \in \}$ = set of labelings of the points in by functions in
- Def. Growth function $() = \max_{|||=} / []/$

Growth function & uniform convergence

- **PAC learning via growth function**: For , > 0 if $|| = \ge 8 / \frac{2}{\ln 2} (\ln 2 (2)) + \ln 1 /)$ then with probability 1 - : $\forall \in H: | () - () | \le 1$
- Thm (Sauer's lemma). If VC-dim(H) = then:

Sauer's Lemma Proof

• Let = -dim() we'll show that if | | = : $| []| \leq (\sum_{\leq i}) = \sum_{i \in I} (i)$ • $\begin{pmatrix} & \\ \leq & \end{pmatrix} = \begin{pmatrix} & -1 \\ & \leq \end{pmatrix} + \begin{pmatrix} & -1 \\ \leq & -1 \end{pmatrix}$ Proof (induction by set size):

$\left| \left[\right] - \left[\left[\left\{ \right\} \right] \right] \le \left(\begin{array}{c} -1 \\ \leq -1 \end{array} \right) \right] \right|$

- If []> [\ { }] then it is because of the sets that differ only on so let's pair them up
- For $h \in []$ containing let $() = h \setminus \{\}$ = $\{h \in []: \in h$ $() \in []\}$
- Note: [] [] <
- What is the VC-dimension of ?
 - $If VC-dim () = 'then \subseteq \setminus \{ \} of 'is shattered$

- All 2 subsets of are 0/1 extendable on

 $- \geq ' + 1 \Rightarrow \text{VC-dim}() \leq -1 \Rightarrow \text{apply induction}$

Examples

- Intervals of the reals:
 - Shatter 2 points, don't shatter $3 \Rightarrow$ -dim = 2
- Pairs of intervals of the reals:
 - Shatter 4 points, don't shatter $5 \Rightarrow$ -dim = 4
- Convex polygons
 - Shatter any points on a circle \Rightarrow -dim = ∞
- Linear separators in dimensions:
 - Shatter + 1 points (unit vectors + origin)
 - Take subset S and set = 0 if \in :

separator ≤ 0

VC-dimension of linear separators

No set of +2 points can be shattered

• Thm (Radon). Any set $\subseteq \mathbb{R}$ with | = +2 can be partitioned into two subsets , s.t.:

Convex() \cap Convex() $\neq \emptyset$

- Form $\times (+2)$ matrix A, columns = points in
- Add extra all-1 row \Rightarrow matrix B
- = $\begin{pmatrix} 1 & 2 & -1 & -1 & -1 & -1 \\ 1 & 2 & -1 & -1 & -1 \end{pmatrix}$, non-zero vector: = 0
- Reordering: $1' 2'' 2'' \ge 0$, +1' 2'' +2 < 0

Radon's Theorem (cont.)



Growth function & uniform convergence

- **PAC learning via growth function**: For , > 0 if $|| = \ge 8 / \frac{2}{\ln 2} (\ln 2 (2)) + \ln 1 /)$ then with probability 1 - : $\forall \in H: | () - () | \le 1$
- Assume event A: $\exists \in H: | () - () | >$



$\left[\right] \geq \Pr[0]] / 2$

- Lem. If $= \Omega(1/2)$ then $\begin{bmatrix} \\ \\ \end{bmatrix} \ge \Pr[\underline{f}_0]/2.$
- Proof: $\begin{bmatrix} \\ \end{bmatrix} \ge \Pr[,] = \Pr[] \Pr[\underline{f}_{0}] /]$
- Suppose occurs: $\exists \in H: | () - () | >$
- When we draw ':

 ()
 ()
 ()
 ()
 ()
 ()

 By Chernoff:
 - $\int_{1}^{1} |z| = \frac{1}{2}$

VC-theorem Proof

- Suffices to show that $\begin{bmatrix} \end{bmatrix} \le /2$
- Consider drawing 2 samples "and then randomly partitioning into 'and
- *: same as for such (',) \Rightarrow Pr [*] = Pr []
- Will show: \forall fixed $"_{,S'} [* / "]$ is small
- Key observation: once "is fixed there are only $\binom{n}{2} \leq \binom{2}{2}$ events to care about
- Suffices: for every fixed $h \in [']$:

VC-theorem Proof (cont.)

- Randomly pair points in "into (,) pairs
- With prob. $\frac{1}{2}$: \rightarrow , \rightarrow 'or \rightarrow ', \rightarrow
- Diff. between () and () for = 1, ...,
- Only changes if mistake on only one of (,)
 - With prob. $\frac{1}{2}$ difference changes by ± 1
- By Chernoff: $Pr\left[\begin{vmatrix} & () - & () \end{vmatrix} > \frac{-\Omega(2)}{4}\right] = -\Omega(2)$ • $-\Omega(2) \le \frac{-\Omega(2)}{2}$ for from the Thm. statement