

機械学習とデータマイニングの基礎 (大阪大学)

PAC 学習と ERM 学習法

Matthew J. Holland

大阪大学 産業科学研究所



目次

1. 今回の趣旨
2. PAC 学習の概要
3. 有限モデルにおける汎化性能
4. 無限モデルにおける汎化性能
5. まとめ

1

目次

1. 今回の趣旨
2. PAC 学習の概要
3. 有限モデルにおける汎化性能
4. 無限モデルにおける汎化性能
5. まとめ

1

今回の趣旨

今回は「学習法の性能評価」を切り口としてある種の「一貫性」の定義から出発.

この「性能」を保証できる代表格の学習法として, ERM を見ていく.

モデルの要素が有限であるか無限であるかで, 証明の技法が大きくことなるため, それぞれの状況での性能保証を探索していく.

特に二値分類の場合, 無限モデルにおける ERM 学習法とモデルの「VC 次元」は深い関係にあり, その基本概念や有名な知見を伝授することを目的とする.

2

目次

1. 今回の趣旨
2. PAC 学習の概要
3. 有限モデルにおける汎化性能
4. 無限モデルにおける汎化性能
5. まとめ

3

理想と現実（作図用）

- ▶ 世の中に流通するデータはノイズで矛盾などもある。
→ すべてのデータで完璧な成績は望めない。
- ▶ どのようなデータが手に入るかも不確定である。
→ 確実にうまく学習する保証は厳しい。

4

学習法の評価（「標本外の汎化」篇）

最低限の性能とは？（直感）

学習データさえ十分にあれば、十分に良い成績をほぼ確実に出してほしい。

定式化

- ▶ n 個の学習データ $\mathbf{Z}_n := (Z_1, \dots, Z_n)$ に基づく $\mathbf{Z}_n \mapsto H_n$ を扱う。
- ▶ 広義のリスク関数 $R(\cdot)$ と比較水準 $R_{\text{con}}^* \leq R(\cdot)$ で評価する。

学習法の「一貫性」(consistency¹)

$$\forall \varepsilon > 0, \lim_{n \rightarrow \infty} \mathbf{P} \{R(H_n) - R_{\text{con}}^* > \varepsilon\} = 0 \quad (1)$$

¹確率変数の期待値を推定する問題における標本平均の弱い一貫性が special case である。

5

学習法の評価（「標本外の汎化」篇）

PAC 学習²

Probably Approximately Correct
 $\geq 1-\delta$ $\leq R_{\text{con}}^* + \varepsilon$ $R(\cdot)$

PAC(n, ε, δ) 条件

$$\mathbf{P} \{R(H_n) - R_{\text{con}}^* > \varepsilon\} \leq \delta \quad (2)$$

- ▶ 標本数 $n \in \mathbb{N}$, 精度 $\varepsilon > 0$, 確度 $1 - \delta$ (ただし $0 < \delta < 1$) とする。
- ▶ 当然, $R(\cdot)$ と R_{con}^* の性質によって, この条件の強度が大きく変わる。
- ▶ この (2) は学習法 $\mathbf{Z}_n \mapsto H_n$ の性質である。

²有名なのは Valiant [10] とその拡張版 [9, 11]. 名称は Angluin [1], Angluin and Laird [2] が初見の模様。

6

学習法の評価（「標本外の汎化」篇）

一貫性と PAC 条件

$$\forall (\varepsilon, \delta), \exists n < \infty, \text{PAC}(n, \varepsilon, \delta) \text{ が成立} \iff (1) \quad (3)$$

標本複雑度 (sample complexity)

$$N_{\varepsilon, \delta}^* := \min \{1 \leq n \leq \infty : \text{PAC}(n, \varepsilon, \delta) \text{ が成立}\} \quad (4)$$

もう少し強い要求³

- ▶ $N_{\varepsilon, \delta}^* = \mathcal{O}(\text{poly}(1/\varepsilon, 1/\delta))$
- ▶ H_n の計算量 $= \mathcal{O}(\text{poly}(1/\varepsilon, 1/\delta))$ (ただし $n \geq N_{\varepsilon, \delta}^*$)

³多項式のレートは必然ではないが、典型的ではある (Kearns and Vazirani [7] を参照)。

目次

1. 今回の趣旨
2. PAC 学習の概要
3. 有限モデルにおける汎化性能
4. 無限モデルにおける汎化性能
5. まとめ

簡単な問題設定から出発する

表記

- ▶ $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ (データの属する集合)
- ▶ $h : \mathcal{X} \rightarrow \mathcal{Y}$ (意思決定の候補)
- ▶ \mathcal{H} (候補の全体集合)
- ▶ $Z = (X, Y) \sim \mu$ (評価用データ)
- ▶ $Z_i = (X_i, Y_i), i \in [n]$ (学習データ)

簡単な問題設定から出発する

事前知識

データの関係性について、以下のことは既知とする。⁴

$$\exists h^* \in \mathcal{H}, \quad Y = h^*(X) \quad (5)$$

この仮定によって明らかなのは、以下の事実である。

$$\mathbf{P} \{h(X) \neq Y\} = \mathbf{P} \{h(X) \neq h^*(X)\} \quad (6)$$

ある種の学習法のアイディアはすぐに思い浮かぶ。

⁴この等式が確率 1 で成り立つという仮定のほうが少し弱い。

簡単な問題設定から出発する

単純な学習法（サル真似）

1. $H_0 \in \mathcal{H}$ を選ぶ.
2. 以降, 各 $t = 1, 2, \dots$ に対して, 以下の手順を行う.
 - ▶ (X_t, Y_t) を得る.
 - ▶ $H_{t-1}(X_t) = Y_t$ なら, $H_t = H_{t-1}$ とする.
 - ▶ $H_{t-1}(X_t) \neq Y_t$ なら, $H_t(X_t) = Y_t$ を満たす $H_t \in \mathcal{H}$ を選ぶ.

最悪時に試す候補の数

$|\mathcal{H}| - 1$ である.

(無限モデル $|\mathcal{H}| = \infty$ なら, 終了しないかもしれない)

11

簡単な問題設定から出発する

統計的な側面（サル真似）1

- ▶ 先述の「サル真似」学習法 ($n+1$ 回の更新) の出力:

$$H_n^{\text{con}} := H_{n+1} \quad (7)$$

- ▶ はっきりとわかっていること:

$$H_n^{\text{con}}(X_i) = Y_i, \quad \forall i \in [n] \quad (8)$$

- ▶ まだわからないこと:

事象「 $\mathbf{P}\{H_n^{\text{con}}(X) \neq Y\} > \varepsilon$ 」の生起確率.

12

簡単な問題設定から出発する

統計的な側面（サル真似）2

ε の水準では「不合格」の候補を集めてみよう.

$$\mathcal{H}_{\text{bad}}(\varepsilon) := \{h \in \mathcal{H} : \mathbf{P}\{h(X) \neq Y\} > \varepsilon\} \quad (9)$$

直感

$H_n^{\text{con}} \in \mathcal{H}_{\text{bad}}(\varepsilon)$ はあり得るが, 全問正解なので可能性が低い (はず).

13

簡単な問題設定から出発する

統計的な側面（サル真似）3

$$\begin{aligned} \mathbf{P}\{H_n^{\text{con}} \in \mathcal{H}_{\text{bad}}(\varepsilon)\} &\leq \mathbf{P}\left[\bigcup_{h \in \mathcal{H}_{\text{bad}}(\varepsilon)} \{h(X_1) = Y_1, \dots, h(X_n) = Y_n\}\right] \\ &\leq \sum_{h \in \mathcal{H}_{\text{bad}}(\varepsilon)} \mathbf{P}\{h(X_1) = Y_1, \dots, h(X_n) = Y_n\} \\ &= \sum_{h \in \mathcal{H}_{\text{bad}}(\varepsilon)} \prod_{i=1}^n \mathbf{P}\{h(X_i) = Y_i\} \\ &= |\mathcal{H}_{\text{bad}}(\varepsilon)| (\mathbf{P}\{h(X) = Y\})^n, \quad h \in \mathcal{H}_{\text{bad}}(\varepsilon) \\ &\leq |\mathcal{H}_{\text{bad}}(\varepsilon)| (1 - \varepsilon)^n \end{aligned}$$

14

簡単な問題設定から出発する

統計的な側面 (サル真似) 4

$$\begin{aligned} \mathbf{P} \{ \mathbf{P} \{ H_n^{\text{con}}(X) \neq Y \} > \varepsilon \} &\leq \mathbf{P} \{ H_n^{\text{con}} \in \mathcal{H}_{\text{bad}}(\varepsilon) \} \leq |\mathcal{H}_{\text{bad}}(\varepsilon)| (1 - \varepsilon)^n \\ &\leq |\mathcal{H}| (1 - \varepsilon)^n \\ &\leq |\mathcal{H}| \exp(-n\varepsilon) \end{aligned}$$

サル真似と PAC 学習

$R(h) = \mathbf{P} \{ h(X) \neq Y \}$, $R_{\text{con}}^* = R(h^*) = 0$ とおけば, 標本複雑度のバウンドを得る.

$$N_{\varepsilon, \delta}^* \leq \frac{\log(|\mathcal{H}|) + \log(1/\delta)}{\varepsilon} \quad (10)$$

この論法の弱点

いつでも「全問正解が可能である」という前提.
これを外すと, 議論が途中で止まってしまう.

15

簡単な問題設定から出発する (作図用)

ノイジーな関係性⁵

$$Y = h^*(X)U, \quad Y = h^*(X) + U, \quad Y = h^*(X + U), \text{ などなど} \cdots$$

X と Y の関係がノイジーであれば, 「矛盾」はあり得るし,
普通のモデルでは表現しきれないことが多々ある.

⁵特に仮定を置かず, $\text{PAC}(n, \varepsilon, \delta)$ で $R_{\text{con}}^* = \inf_{h \in \mathcal{H}} R(h)$ とする設定を Agnostic PAC という.

16

ERM 学習法

サル真似の特徴

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1} \{ H_n^{\text{con}}(X_i) \neq Y_i \} = 0 \quad (11)$$

その一般化 ($\ell(\cdot)$ は自由)

$$H_n^{\text{ERM}} \in \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h; Z_i) \quad (12)$$

ERM 学習法⁶

$$\underbrace{\text{Empirical Risk Minimization}}_{(Z_1, \dots, Z_n) \quad \mathbf{E}[L] \quad \min_{h \in \mathcal{H}}}$$

⁶統計的学習理論の中枢にあたる学習法. Vapnik and Chervonenkis [12] は画期的である.

17

ERM 学習法

論法を点検する

ここで「全問正解」を「ERM 解」に緩和してみよう.

元の議論の途中にあった事象を以下のように置き換えざるを得ない.

$$\begin{aligned} \{h(X_1) = Y_1, \dots, h(X_n) = Y_n\} \\ \downarrow \\ \{h(X_1) = H_n^{\text{ERM}}(X_1), \dots, h(X_n) = H_n^{\text{ERM}}(X_n)\} \end{aligned}$$

$H_n^{\text{ERM}}(\cdot)$ が全学習データに依存するため, 個別の事象の独立性が失われてしまう.

18

ERM 学習法

表記

$$L(h) := \ell(h; Z), \quad R(h) = \mathbf{E}_\mu L(h), \quad R_n(h) := \frac{1}{n} \sum_{i=1}^n \ell(h; Z_i) \quad (13)$$

ひとつのアイデア

Z_1, \dots, Z_n と Z が独立同分布ならば、以下の論法はできそう。

- ▶ ERM 学習法で $R_n(\cdot)$ を小さくする。
- ▶ 各 $h \in \mathcal{H}$ に対して、 $R_n(h) \approx R(h)$ を示す。
- ▶ ERM 学習法では $R(\cdot)$ も小さくなることを示す。

肝心の近似の良し悪しについて、もう少し突き詰めていこう。

19

ERM 学習法

集中不等式 (Chebyshev 型)⁷

$$\mathbf{P}\{|R_n(h) - R(h)| > \varepsilon\} \leq \frac{\mathbf{E}(R_n(h) - R(h))^2}{\varepsilon^2} = \frac{\text{var}_\mu L(h)}{\varepsilon^2 n} \quad (14)$$

集中不等式 (Hoeffding 型)⁸

$\mathbf{P}\{a \leq L(h) \leq b\} = 1$ が成り立つ場合、以下のことが言える。

$$\mathbf{P}\{|R_n(h) - R(h)| > \varepsilon\} \leq 2 \exp\left(\frac{-2n\varepsilon^2}{(b-a)^2}\right) \quad (15)$$

⁷Markov の不等式から容易に導き出すことができる。

⁸Hoeffding [6] の技法による。Boucheron et al. [4] には明快な証明が収録されている。

20

ERM 学習法

議論の流れ⁹

例として先ほどの (15) を借用すると、以下の不等式が出せる。

$$\begin{aligned} \mathbf{P}\{|R_n(H_n^{\text{ERM}}) - R(H_n^{\text{ERM}})| > \varepsilon\} &\leq \mathbf{P}\left\{\sup_{h \in \mathcal{H}} |R_n(h) - R(h)| > \varepsilon\right\} \\ &= \mathbf{P}\left[\bigcup_{h \in \mathcal{H}} \{|R_n(h) - R(h)| > \varepsilon\}\right] \\ &\leq \sum_{h \in \mathcal{H}} \mathbf{P}\{|R_n(h) - R(h)| > \varepsilon\} \\ &\leq 2|\mathcal{H}| \exp\left(\frac{-2n\varepsilon^2}{(b-a)^2}\right) \end{aligned}$$

⁹ここで、ERM の性質を何も使っていないので、任意の学習法にも通じる議論。

21

ERM 学習法

ERM 全般についてわかっていること

確率 $1 - \delta$ 以上で成り立つリスクバウンド。

$$R(H_n^{\text{ERM}}) \leq R_n(H_n^{\text{ERM}}) + (b-a) \sqrt{\frac{1}{2n} \log\left(\frac{2|\mathcal{H}|}{\delta}\right)} \quad (16)$$

更には、以下の便利な不等式が使える。

$$\begin{aligned} R(H_n^{\text{ERM}}) - R_{\mathcal{H}}^* &= R(H_n^{\text{ERM}}) - R_n(H_n^{\text{ERM}}) + R_n(H_n^{\text{ERM}}) - R_{\mathcal{H}}^* \\ &\leq 2 \sup_{h \in \mathcal{H}} |R_n(h) - R(h)| \end{aligned} \quad (17)$$

結果として、以下の不等式も導き出せる (これも確率 $1 - \delta$ 以上)。

$$R(H_n^{\text{ERM}}) - R_{\mathcal{H}}^* \leq 2(b-a) \sqrt{\frac{1}{2n} \log\left(\frac{2|\mathcal{H}|}{\delta}\right)} \quad (18)$$

22

ERM 学習法

ノイズ有りの任意 ERM 学習法 ((18) の再掲)

$$R(H_n^{\text{ERM}}) - R_H^* \leq 2(b-a) \sqrt{\frac{1}{2n} \log \left(\frac{2|\mathcal{H}|}{\delta} \right)} \quad (\text{確率 } 1 - \delta \text{ 以上})$$

ノイズ無しの「サル真似」と比べて (ゼロイチ損失)

$$R(H_n^{\text{con}}) \leq \frac{1}{n} \log \left(\frac{|\mathcal{H}|}{\delta} \right) \quad (19)$$

(確率 $1 - \delta$ 以上)

23

目次

1. 今回の趣旨
2. PAC 学習の概要
3. 有限モデルにおける汎化性能
4. 無限モデルにおける汎化性能
5. まとめ

24

無限モデルにおける汎化性能

この論法の限界

$|\mathcal{H}| = \infty$ の場合, 先述の (16) と (19) はいずれも破綻する (保証として無意味).

新たな切り口

この \mathcal{H} の表現力を単なる「候補の数」とするのはいささか乱暴なので, 表現力の指標をもう少し工夫してから議論を再建してみよう.

25

モデル表現力の新しい指標

問題設定を少し狭める

$\mathcal{Y} = \{-1, +1\}$ とおく. つまり二値分類に限定する.

直感

二値ラベルを割り当てる「自由度」が高いほど, 表現力が高い.

26

モデル表現力の新しい指標

識別能力の指標

ここで n 個の入力データを $\mathbf{x}_n := (x_1, \dots, x_n)$ と記す.

$$\text{shatter}(n; \mathcal{H}, \mathcal{X}) := \max_{\mathbf{x}_n \in \mathcal{X}^n} |\underbrace{\{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\}}_{\text{dichotomies of } \mathcal{H}}| \quad (20)$$

要点

- ▶ \mathcal{H} によって割り当てられる二値ラベルの組み合わせの数.
- ▶ $1 \leq \text{shatter}(n) \leq 2^n$ は常に成り立つ.
- ▶ ある m で $\text{shatter}(m) < 2^m$ なら, 任意の $n > m$ でも $\text{shatter}(n) < 2^n$.
- ▶ 名称として shatter coefficient や growth function などと呼ばれる.

27

モデル表現力の新しい指標 (作図用)

例: 実数直線状の有限区間

$$\text{shatter}(3) < 2^3$$

28

モデル表現力の新しい指標 (作図用)

例: 二次元平面における識別線

$$\text{shatter}(4) < 2^4$$

29

モデル表現力の新しい指標 (作図用)

例: \mathbb{R}^d における長方形

$$\text{shatter}(2d+1) < 2^{2d+1}$$

30

モデル表現力の新しい指標

VC 次元の定義¹⁰

$$VC(\mathcal{H}; \mathcal{X}) := \max \{1 \leq n \leq \infty : \text{shatter}(n; \mathcal{H}, \mathcal{X}) = 2^n\} \quad (21)$$

先ほどの例

- ▶ 実数直線状の有限区間: $VC = 2$
- ▶ 二次元平面における識別線: $VC = 3$
- ▶ \mathbb{R}^d における長方形: $VC = 2d$

¹⁰Vapnik-Chervonenkis の研究成果に由来 [12]. 入門書として Devroye et al. [5, Ch. 12–14] は傑作.

モデル表現力の新しい指標

面白い性質¹¹

任意の $n > 2 VC(\mathcal{H}; \mathcal{X})$ に対して, 以下が成り立つ.

$$\text{shatter}(n; \mathcal{H}, \mathcal{X}) \leq \left(\frac{en}{VC(\mathcal{H}; \mathcal{X})} \right)^{VC(\mathcal{H}; \mathcal{X})} \quad (22)$$

また, $\text{shatter}(n) \leq n^{VC} + 1$ は常に成り立つ.

つまり, $\text{shatter}(\cdot)$ は「指数増加」と「多項式増加」のどちらかしかあり得ない.

¹¹Devroye et al. [5, Thm. 13.3] を参照. 1970 年代までに様々な分野でほぼ同時に明らかになった.

モデル表現力の新しい指標

無限モデルの性能保証にどう役立つ? (直感)

$R_n(h) \approx R(h)$ をナイーブに全 $h \in \mathcal{H}$ で抑えると $|\mathcal{H}|$ 依存は不可避

↓

\mathcal{H} が無限でも, 二値分類の仕方は有限であるため,
一致する候補同士は同一視しても良い

↓

$|\mathcal{H}|$ を $\text{shatter}(n; \mathcal{H}, \mathcal{X})$ に置き換える

モデル表現力の新しい指標

無限モデルの性能保証にどう役立つ? ¹²

ゼロイチ損失の下, 以下のことが言える (証明はすべて割愛).

$$\mathbf{P} \left\{ \sup_{h \in \mathcal{H}} |R_n(h) - R(h)| > \varepsilon \right\} \leq 8 \text{shatter}(n; \mathcal{H}, \mathbb{R}^d) \exp(-n\varepsilon^2/32) \quad (23)$$

さらに, (17) を使うと, 以下のこともわかる.

$$\mathbf{P} \left\{ R(H_n^{\text{ERM}}) - R_{\mathcal{H}}^* > \varepsilon \right\} \leq 8 \text{shatter}(n; \mathcal{H}, \mathbb{R}^d) \exp(-n\varepsilon^2/128) \quad (24)$$

¹²Devroye et al. [5, Thm. 12.6] を参照.

モデル表現力の新しい指標

VC 次元と ERM の性能

VC 次元が有限である \implies 無限モデルでも ERM 学習法の性能保証が得られる

ERM 学習法の平均的な性能¹³

$$\mathbf{E} \left[R(H_n^{\text{ERM}}) - R_H^* \right] \leq 16 \sqrt{\frac{\log(8e) + \log(\text{shatter}(n; \mathcal{H}, \mathbb{R}^d))}{n}} \quad (25)$$

¹³Devroye et al. [5, Cor. 12.1] を参照.

モデル表現力の新しい指標

ERM 学習法の標本複雑度¹⁴

$$N_{\varepsilon, \delta}^* \leq \max \left\{ \mathcal{O} \left(\frac{\text{VC}}{\varepsilon^2} \log \left(\frac{\text{VC}}{\varepsilon^2} \right) \right), \mathcal{O} \left(\frac{1}{\varepsilon^2} \log \left(\frac{1}{\delta} \right) \right) \right\} \quad (26)$$

ノイズ無しの「サル真似」と比べて

$$N_{\varepsilon, \delta}^* \leq \frac{\log(|\mathcal{H}|) + \log(1/\delta)}{\varepsilon} \quad (27)$$

¹⁴Devroye et al. [5, Cor. 12.3] を参照.

PAC 学習可能性

学習可能性¹⁵

任意の分布 μ の下, $N_{\varepsilon, \delta}^* < \infty$ を満たすような学習法が存在すれば, その \mathcal{H} を「学習可能 (learnable)」という.

VC 次元との関係 (二値分類の場合) 1

先述の (26) からすぐにわかるのは,

\mathcal{H} の VC 次元が有限である $\implies \mathcal{H}$ は学習可能である

ということである.

¹⁵大いに研究されてきた. 成熟した理論を Shalev-Shwartz et al. [8] はうまくまとめている.

PAC 学習可能性

VC 次元との関係 (二値分類の場合) 2

さらに面白いのは, その逆である. 任意の学習法に対して,

$$\text{VC}(\mathcal{H}; \mathcal{X}) = \infty \implies N_{\varepsilon, \delta}^* = \infty \text{ となる } \mu \text{ が存在する.}^{16}$$

つまり「学習可能性は VC 次元の有限性によって完全に特徴づけられる」.

(これは「**学習理論の基本定理**」と呼ばれることがある)

¹⁶Anthony and Bartlett [3, Thm. 5.2] を参照.

目次

1. 今回の趣旨
2. PAC 学習の概要
3. 有限モデルにおける汎化性能
4. 無限モデルにおける汎化性能
5. まとめ

39

まとめ

- ▶ 推定量の「一致性」という概念は学習法にも拡張できる。
- ▶ 特にデータや計算のコストを考慮した一致性の探求は PAC 学習。
- ▶ 従来の学習理論は $R(h) = E_{\mu} L(h)$ が基軸となっている。
- ▶ 計算量を捨象し、ERM 学習法の統計的な性質（リスク上界を求めるなど）は統計的学習理論の中心トピック。
- ▶ 有限モデルと無限モデルにおける ERM の解析は手段が大きく異なるが、基本的な議論のアイデアは共通している。
- ▶ 有限モデルでも無限モデルでも、 R_n をめぐる「集中不等式」が登場する。
- ▶ VC 次元は二値分類問題の「解きやすさ」を特徴づけている。

40

参考文献

- [1] Angluin, D. (1988). Queries and concept learning. *Machine Learning*, 2(4):319–342.
- [2] Angluin, D. and Laird, P. (1988). Learning from noisy examples. *Machine Learning*, 2(4):343–370.
- [3] Anthony, M. and Bartlett, P. L. (1999). *Neural Network Learning: Theoretical Foundations*. Cambridge University Press.
- [4] Boucheron, S., Lugosi, G., and Massart, P. (2013). *Concentration inequalities: a nonasymptotic theory of independence*. Oxford University Press.
- [5] Devroye, L., Györfi, L., and Lugosi, G. (1996). *A Probabilistic Theory of Pattern Recognition*. Springer.
- [6] Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30.
- [7] Kearns, M. J. and Vazirani, U. V. (1994). *An Introduction to Computational Learning Theory*. MIT Press.
- [8] Shalev-Shwartz, S., Shamir, O., Srebro, N., and Sridharan, K. (2010). Learnability, stability and uniform convergence. *Journal of Machine Learning Research*, 11:2635–2670.
- [9] Valiant, L. G. (1984a). Deductive learning. *Philosophical Transactions of the Royal Society of London, Series A*, 312(1522):441–446.
- [10] Valiant, L. G. (1984b). A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142.
- [11] Valiant, L. G. (1985). Learning disjunction of conjunctions. In *Proceedings of the 9th International Joint Conference on Artificial Intelligence*, pages 560–566. Morgan Kaufmann.
- [12] Vapnik, V. N. and Chervonenkis, A. Y. (1971). On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability & Its Applications*, 16(2):264–280.

41