

# למידה חישובית

אלי דיין<sup>1</sup>

### תקציר

מסמך זה יביא את סיכומי השיעורים מהקורס למידה חישובית, שהועבר על ידי פרופ' ישי מנצור בסמסטר א' בשנה"ל תשע"ג.

# תוכן עניינים

5	מה זה למידה חישובית?	1
5	סוגי הבעיות	1.1
5	סוגים של ML	1.2
5	Supervised vs. Unsupervised	1.2.1
6	Active vs. Passive	1.2.2
6	Teacher	1.2.3
6	Batch vs. Online	1.2.4
6	בניית מודל ML	1.3
7	Loss Model	1.3.1
7	0-1Loss	1.3.1.1
7	Quadratic Loss	1.3.1.2
7	Logarithmic Loss	1.3.1.3
8	הנחות על השערות	1.3.2
8	שיטות הסקה	1.3.3
8	הסקה בייסיאנית (Bayes)	1.3.3.1
9	מודל PAC	1.3.3.2
9	מודלים מכוונים (Online)	1.3.3.3
10	מבנה הקורס	1.4
11	Bayesian Inference	2
11	מבוא	2.1
11	כלל Bayes	2.2
11	דוגמה: זיהוי סרטן	2.3
13	דוגמה: התפלגות נורמלית	2.4
13	התפלגות נורמלית	2.4.1
13	תיאור הבעיה	2.4.2
13	שיטת ML	2.4.3
14	שיטת MAP	2.4.4
15	שיטת Posterior Bayes	2.4.5
16	Learning a Concept Class	2.5
17	דוגמה: Biased Coins	2.6
17	חוק Laplace	2.6.1
18	פונקציות Loss	2.6.2
20	Naïve Bayes	2.7
20	סיווג בייסיאני: מרחב בינארי	2.7.1
21	פענוח של Naïve Bayes	2.7.2
21	התפלגות נורמלית	2.7.3

<b>23</b>	<b>מודל ה-PAC</b>	<b>3</b>
23	מודל ה-PAC	3.1
23	דוגמה אינטואיטיבית	3.2
24	מציאת השערה טובה	3.2.1
24	אופן הלימוד	3.2.2
24	מספר הדגימות	3.2.3
26	הצגה פורמלית של מודל ה-PAC	3.3
26	הקדמה	3.3.1
26	הגדרת מודל ה-PAC	3.3.2
27	מחלקות השערות סופיות	3.4
27	המקרה $c_t \in H$	3.4.1
27	המקרה $c_t \notin H$	3.4.2
29	דוגמה - למידת Boolean Disjunctions	3.4.3
29	דוגמה - למידת Parity	3.4.4
30	Occam Razor	3.5
30	אלגוריתמי Occam ומודל ה-PAC	3.5.1
30	דוגמה - למידת OR של $k$ משתנים	3.5.2
<b>32</b>	<b>מודל Online</b>	<b>4</b>
32	למידה של מפריד לינארי	4.1
32	אלגוריתם Perceptron	4.1.1
34	אלגוריתם Margin Perceptron	4.1.2
35	מודל Margin Bound	4.2
36	האלגוריתם (CON) Consistent	4.2.1
36	אלגוריתם חציה (HAL)	4.2.2
36	הקשר בין Mistake Bound ומודל ה-PAC	4.3
37	למידה של OR	4.4
37	אלגוריתם Winnow	4.5
<b>40</b>	<b>Regret Minimization</b>	<b>5</b>
40	מבוא	5.1
40	המודל האלגוריתמי	5.1.1
40	External Regret	5.1.2
41	אלגוריתמים	5.2
41	אלגוריתם (G) Deterministic Greedy	5.2.1
42	אלגוריתם (GR) Randomized Greedy	5.2.2
43	אלגוריתם (RWM) Randomized Weighted Majority	5.2.3
45	חסמים תחתונים לאלגוריתמי Online ממושקלים	5.3
45	טווח קצר - $T = \frac{1}{2} \cdot \log  H $	5.3.1
46	כתלות בזמן - $ H  = 2$	5.3.2
46	Multi-Arm Bandit	5.4
47	אלגוריתם Test & Play	5.4.1
48	אלגוריתם (UCB) Upper Confidence Bound	5.4.2
<b>50</b>	<b>Boosting</b>	<b>6</b>
50	למידה חלשה וחזקה	6.1
50	שיפור בפרמטר הבטחון	6.1.1
52	שיפור בפרמטר הדיוק	6.1.2
53	בנייה רקורסיבית	6.2
55	אלגוריתם AdaBoost	6.3

<b>60</b>	<b>Nearest Neighbor</b>	<b>7</b>
60	הקדמה	7.1
60	שיטות כלליות	7.2
60	מודלים ל-Nearest Neighbor	7.3
61	0-1 Loss	7.3.1
61	Bayes Risk	7.3.2
61	מקרה פשוט	7.3.2.1
62	המקרה הכללי	7.3.2.2
63	שיטת $k$ שכנים קרובים ( $k$ -NN)	7.4
63	מקרה פשוט	7.4.1
64	מקרה כללי	7.4.2
64	מדידת המרחק	7.5
64	Locality Sensitive Hashing	7.6
65	שלב א' - Amplification	7.6.1
66	שלב ב'	7.6.2
66	האלגוריתם	7.6.3
<b>68</b>	<b>VC Dimension</b>	<b>8</b>
68	מודל PAC (חזרה)	8.1
69	מימד VC	8.2
69	מוטיבציה	8.2.1
69	הגדרות	8.2.2
70	דוגמאות	8.2.3
70	סיפא על קו	8.2.3.1
70	מפריד לינארי במישור	8.2.3.2
73	מלבנים מקבילים לצירים	8.2.3.3
74	מספר סופי של אינטרוולים	8.2.3.4
74	פוליגון קונבקסי במישור	8.2.3.5
74	חסם תחתון על גודל הדגימה	8.2.4
75	עוד דוגמאות	8.2.5
75	פונקציית Parity	8.2.5.1
76	OR של ליטרלים	8.2.5.2
77	מפריד לינארי במימד $n$	8.2.5.3
<b>79</b>	<b>מימד VC (המשך)</b>	<b>9</b>
79	חזרה	9.1
80	חסמים לגודל הדגימה	9.2
82	סיבוכיות Radamacker	9.3
82	ממוצעי Radamacker	9.3.1
83	אי שוויון McDiarmid	9.3.2
84	סיבוכיות Radamacker	9.3.3
<b>86</b>	<b>שיעור 10</b>	<b>10</b>
<b>87</b>	<b>שיעור 11</b>	<b>11</b>

<b>88</b>	<b>12 רגרסיה</b>
88	12.1 הקדמה
89	12.2 רגרסיה לינארית
90	12.2.1 חוסר יציבות הפתרון
90	12.3 רגולריזציה
90	12.3.1 Ridge Regression
91	12.3.2 Lasso Regression
91	12.3.3 חסם הכללה ל-Ridge Regression (או משהו שדומה לו)
92	12.3.4 נקודת מבט בייסיאנית
93	12.4 Logistic Regression
<b>95</b>	<b>13 Model Selection</b>
95	13.1 הקדמה
96	13.1.1 דוגמה
96	13.1.2 המודל
97	13.2 Structural Risk Minimization
99	13.3 Cross Validation ( $CV$ )
100	13.4 Minimum Description Length ( $MDL$ )
101	13.4.1 בעזרת MAP

# פרק 1

## מה זה למידה חישובית?<sup>1</sup>

התחום צמח מתוך תחום הבינה המלאכותית ודומה לסטטיסטיקה, מבחינת השאלות ששואלים.

### 1.1 סוגי הבעיות

1. סיווג - Classification.

נניח שלכל קלט יש סיווג נכון (ואולי יחיד). נרצה לדעת מהו. לדוגמה:

(א) spam - סינון דואר זבל. כל הודעה יכולה להיות הודעה אמיתית או הודעת דואר זבל. נרצה לסווג את ההודעה בהתאם.

(ב) סרטן (או מציאת מחלות באופן כללי).

(ג) כרטיס אשראי - האם טרנסאקציה היא לגיטימית? לפסול אותה או לאשר אותה?

(ד) דוגמה למקרה בו לקלט מספר סיווגים שונים: נושא המאמר (למשל פוליטיקה, ספורט ועוד).

2. בעיות בקרה.

צריך לקבל החלטה שמשפיעה על מה שיהיה הקלט בפעם הבאה. דוגמאות:

(א) לשחק משחקים (למשל ארבע בשורה, או שח-מט).

(ב) רובוטים - למשל שליטה על מסוק.

3. מערכות המלצה - Collaborative Filter (כמו המלצות לחברים ב-Facebook).

בקורס נתמקד בבעיות סיווג.

### 1.2 סוגים של ML<sup>2</sup>

#### 1.2.1 Supervised vs. Unsupervised

Supervised מקבלים דוגמאות עם סיווג  $y$ ,  $(a_1, a_2, \dots, a_n)$ .

Unsupervised אין סיווג. המטרה היא לחלק את הנתונים לקבוצות בעלות דמיון.

בקורס נתמקד בבעיות Supervised.

<sup>1</sup>שיעור שהתקיים בתאריך 21.10.2012.  
<sup>2</sup>קיצור של Machine Learning.

**Active vs. Passive 1.2.2**

Passive יש מאגר מידע שמשמשים בו בצורה פאסיבית. למשל: data base -> השערה <- חיזוי

Active האלגוריתם בוחר את הקלטים שיסווגו.

בקורס נתמקד באלגוריתמים שהם Passive.

**Teacher 1.2.3**

Teacher למידה על ידי מומחה.

בקורס לא נתעסק בזה.

**Batch vs. Online 1.2.4**

Batch כל המידע נתון מראש.

Online המידע אינו נתון מראש. האלגוריתם רואה דוגמה, מנחש סיווג, מקבל את הסיווג הנכון, וממשיך לדוגמה הבאה.

בקורס נדבר על שתי הגישות.

**1.3 בניית מודל ML**

בכל מודל, נרצה להגדיר את המאפיינים הבאים:

1. איך הדוגמאות מיוצרות? למעשה, אנו רוצים כאן הנחה לגבי הסביבה.

2. להשוות סיווגים של אלגוריתמים שונים. הרבה שגיאות קטנות לעומת מעט שגיאות גדולות.

3. מיפוי מדוגמה לסיווג (השערה).

בהינתן שתי השערות, + ו-, וההתפלגויות  $D_+$  ו-  $D_-$  של + ו- בהתאמה, נסמן:

$$\lambda = \Pr[+] = 1 - \Pr[-]$$

כמו כן:

$$D(x) = \lambda \cdot D_+(x) + (1 - \lambda) \cdot D_-(x)$$

בהינתן נקודה  $(x, y)$ , היינו רוצים לענות על השאלה: מה ההסתברות שההשערה + נכונה בהינתן  $(x, y)$ ? כלומר, היינו רוצים למצוא את  $\Pr[+ | (x, y)]$ . נפתח את הנוסחה:

$$\begin{aligned} p = \Pr[+ | (x, y)] &= \frac{\Pr[(x, y) | +] \cdot \Pr[+]}{\Pr[(x, y)]} \\ &= \frac{D_+(x, y) \cdot \lambda}{D(x, y)} \end{aligned}$$

כאן, אנחנו צריכים לשאול את עצמנו מה המטרה של הסיווג. למשל, אם אנחנו מאתרים סרטן על פי הנתונים של החולה, אז יש חשיבות רבה לסיווג שניתן. אם נחליט שהאיש חולה, נפנה אותו לעוד בדיקות. לעומת זאת, אם נחליט שהאיש בריא, נשחרר אותו. לכן, היינו מעדיפים שבמקרים של ספק נשלח את האיש לבדיקות נוספות.



**Loss Model 1.3.1****0-1Loss 1.3.1.1**

במקרה של סיווג דו-ערכי:

$$\ell(0, 1) = 1 = \ell(1, 0)$$

$$\ell(0, 0) = 0 = \ell(1, 1)$$

במקרה הכללי:

$$\ell(y_1, y_2) = \begin{cases} 0 & y_1 = y_2 \\ 1 & y_1 \neq y_2 \end{cases}$$

כלומר, במידה וצדקנו, ההפסד (Loss) שלנו הוא 0. בכל מקרה אחר, 1. בדוגמה לעיל (עם הסרטן):

$$\begin{aligned} \Pr[+ | (x, y)] &\geq \Pr[- | (x, y)] \\ p &\geq 1 - p \\ p &\geq 1/2 \end{aligned}$$

**Quadratic Loss 1.3.1.2**

המודל אמור להביא קירוב טוב ביותר לכל המקרים, מבחינת מרחק. אם ההסתברות הנכונה היא  $p$ , ההפסד עבור ההסתברות  $q$  הוא:

$$\begin{aligned} \ell(q) &= p(1-q)^2 + (1-p)q^2 \\ \frac{d}{dq}\ell &= -2p(1-q) + 2(1-p)q = 0 \\ -p + pq + q - pq &= 0 \\ p &= q \end{aligned}$$

כלומר, התשובה הטובה ביותר היא  $p$  (כמצופה).

**Logarithmic Loss 1.3.1.3**

אם הסיווג הוא  $+$ , אז השגיאה תהיה  $-\log q$ . אם הסיווג הוא  $-$ , אז השגיאה תהיה  $-\log(1-q)$ .

אם ההסתברות הנכונה היא  $p$ , ההפסד עבור ההסתברות  $q$  הוא:

$$\begin{aligned} \ell(q) &= -p \log q - (1-p) \log(1-q) \\ \frac{d}{dq}\ell &= -\frac{p}{q} + \frac{1-p}{1-q} = 0 \\ p(1-q) - q(1-p) &= 0 \\ p &= q \end{aligned}$$

כלומר, התשובה הטובה ביותר היא  $p$  (כמצופה).

**הערה** בכל השגיאות קיבלנו כי התשובה האופטימאלית היא  $p$ . אז למה יש סוגים שונים של שגיאות? מכיוון שהשגיאות עוסקות במקרים התת-אופטימליים, והם אילו שמעניינים אותנו.

### 1.3.2 הנחות על השערות

הנחה מפורשת:

$$f(x) = \sum_{i=1}^d \alpha_i x_i$$

או באופן מדויק יותר,  $f(x) + \text{Noise}$ . מחלקות השערות:

$$H = \left\{ \sum_{i=1}^d \alpha_i x_i \geq \theta \right\}$$

כאשר  $\alpha_i \geq 0$  ו- $\theta \in [0, 1]$ ,  $\sum \alpha_i = 1$

### 1.3.3 שיטות הסקה

#### 1.3.3.1 הסקה בייסיאנית (Bayes)

ניעזר בהתפלגות prior (המפלה את ה"אמונות" שלנו). למשל, אם הקלט הוא מדגם  $\{x_i, b_i\}$ ,  $S = \{x_i, b_i\}$ ,  $x_i \in \mathbb{R}^d$ ,  $b_i \in \{0, 1\}$ , ואנחנו מחפשים את  $f(x)$  אזי:

$$\Pr[f(x) = 1 | S, x] = \sum_{h \in H} h(x) \cdot \Pr[f = h | S]$$

$$\Pr[f = h | S] = \frac{\Pr[S | f = h] \cdot \Pr[f = h]}{\Pr[S]}$$

אותנו  $\Pr[S]$  לא מעניין, כי זה הנתון שלנו, ולכן אנחנו יכולים להניח כי  $\Pr[S] = 1$ . לכן:

$$\Pr[f = h | S] = \Pr[S | f = h] \cdot \Pr[f = h]$$

כעת, ישנן שתי שיטות לבחירת  $h$ :

1. Maximum Likelihood:

$$h_{\text{ML}} = \arg \max_{h \in H} \Pr[S | f = h]$$

2. Maximum A Posteriori:

$$h_{\text{MAP}} = \arg \max_{h \in H} \Pr[S | f = h] \cdot \Pr[f = h]$$

## 1.3.3.2 מודל PAC

ההתפלגות  $D$  לא ידועה. כמוה גם פונקציית המטרה  $f$  אינה ידועה.  $H$  היא מחלקת ההשערות.

המטרה שלנו היא למצוא  $h \in H$  שתביא למינימום את  $\Pr_D [h(x) \neq f(x)] = \varepsilon(h)$ .

**בעיה** אין לנו את  $D$ , אלא רק מדגם  $S$  מתוך  $D$ . נגדיר:

$$\hat{\varepsilon}(h) = \frac{1}{|S|} \cdot \sum_{x \in S} I[h(x) \neq f(x)]$$

כאשר  $I[\delta]$  הוא האינדיקטור של  $\delta$ , כלומר:

$$I[\delta] = \begin{cases} 1 & \delta \\ 0 & -\delta \end{cases}$$

כעת, אנחנו יכולים למצוא  $h \in H$  שיביא למינימום את  $\hat{\varepsilon}(h)$ .

**שאלה** מה גודל  $|\varepsilon(h) - \hat{\varepsilon}(h)|$ ? בהמשך נבטיח גם חסם הכללה:

$$\forall h \in H. |\varepsilon(h) - \hat{\varepsilon}(h)| \leq \delta$$

המתודולוגיה הכללית שלנו תהיה:

1. נבחר  $h$  שיביא למינימום את  $\hat{\varepsilon}(h)$  (ERM).
2. נוכיח חסם הכללה.

## דוגמאות למחלקות השערות

1. קו ישר על מישור (או  $\mathbb{R}^{n-1}$  שמפריד בתוך  $\mathbb{R}^n$ ).
2. עץ החלטה.

**דוגמה<sup>3</sup>** נניח כי  $X = \{0, 1\}^d$ ,  $H = \{f_{x_1}, f_{x_2}, \dots, f_{x_d}\}$  כאשר  $f_{x_i}(\langle x_1, x_2, \dots, x_n \rangle) = x_i$  ו- $D \sim \text{Unif}(H)$ . נגדיל פונקציה  $f(x)$  כלשהי. מה יקרה אם  $m < \frac{1}{2} \log d$  דוגמאות? ההסתברות ש- $x_1$  מסווג  $m$  דוגמאות נכון:  $\frac{1}{2^m}$ . ההסתברות שקיים  $x_i$  שמסווג את כל הדוגמאות נכון:  $1 - (1 - \frac{1}{2^m})^d$ . אם נשאיף  $d \rightarrow m$ , נקבל  $\frac{1}{e}$ . בהסתברות קבועה, קיים  $x_i$  שמסווג את כל הדוגמאות נכון.

למעשה, היינו רוצים לחשב את  $\min_h \hat{\varepsilon}(h) + \text{Complexity}(h)$ , מכיוון שכאשר הסיבוכיות גדולה מדי, אנחנו נקבל Overfitting למדגם, והשגיאה בעולם שמעבר למדגם תהיה גדולה מדי.

## 1.3.3.3 מודלים מכוונים (Online)

בזמן  $t$  רואים קלט  $x_t$ , מוציאים סיווג  $y_t$ , רואים את הסיווג האמיתי  $f(x_t)$ , צוברים הפסד  $\ell(y_t, f(x_t))$  וממשיכים.

<sup>3</sup> כל הקטע הזה לא ברור בכלל.

**המטרה** להביא למינימום את  $\mathcal{L} = \sum_t \ell(y_t, f(x_t))$ .

**פתרון** נניח כי  $h^* \in H$ , שעבורה  $\mathcal{L}$  קטן. כלומר:

$$\mathcal{L}(h^*) = \sum_t \ell(h^*(x_t), f(x_t))$$

נרצה:

$$\mathcal{L}_{\min} \leq \mathcal{L}(h^*) + R$$

## 1.4 מבנה הקורס

- מודלים בסיסיים:

- .Bayes

- .PAC

- .Regression

- אלגוריתמי Online:

- .Perception

- .Regret

- .Boosting

- Generalization Bounds:

- .VC-dim

- .Radamacher

- אלגוריתמים:

- .Decission Trees

- .SVM

- .Fourier Transform

## פרק 2

# Bayesian Inference<sup>1</sup>

### 2.1 מבוא

בשיעור זה נציג את מודל ההסקה הבייסיאני (Bayes Inference). המודל משתמש בהתפלגות prior, שמשקפת את האמונות שלנו לגבי ההשערה הנכונה או הסיווג הנכון, ובסוף מקבלים התפלגות posterior, שמשקפת את מה שהתהליך למד. יש 3 דרכים להסקה בייסיאנית:

- ML (או Maximum Likelihood).
- MAP (או Maximum A Posteriori).
- Bayes Posterior Rule.

### 2.2 כלל Bayes

$$\Pr[A | B] = \frac{\Pr[B | A] \cdot \Pr[A]}{\Pr[B]} \quad (2.2.1)$$

איך אנחנו נשתמש בזה? אנחנו נרצה לדעת את ההסתברות שההשערה  $h$  היא נכונה בהינתן המדגם  $data$ , כלומר את  $\Pr[h | data]$ . לפי חוק Bayes:

$$\Pr[h | data] = \frac{\Pr[data | h] \cdot \Pr[h]}{\Pr[data]}$$

נשים לב ש- $\Pr[data]$  הינו בחזקת קבוע, מכיוון שלרוב אין לנו מידע על בחירת המדגמים, ולכן בהרבה מקומות ניתן להשמיט אותו.

### 2.3 דוגמה: זיהוי סרטן

נתונה ערכה לזיהוי של מחלת הסרטן. בהינתן מטופל, הערכה עשויה להחזיר אחת משתי תשובות: 0 או 1, המייצגות מטופל בריא או חולה בהתאמה. כמו כן, נתון כי:

---

<sup>1</sup>שיעור שהתקיים בתאריך 28.10.2012.

- אם המטופל חולה בסרטן, בסיכוי של 98% הערכה תחזיר עבורו 1.
  - אם המטופל בריא, בסיכוי של 97% הערכה תחזיר עבורו 0.
  - ההסתברות הכללית לסרטן בקרב האוכלוסייה היא 1%.
- היינו רוצים למצוא את הסיכוי שמטופל חולה בסרטן אם הערכה טוענת שהוא חולה, כלומר את  $\Pr[\text{Cancer} | 1]$ .  
נפרמל את הנתונים שלנו:

$$\Pr[1 | \text{Cancer}] = 98\% = 0.98$$

$$\Pr[0 | \neg\text{Cancer}] = 97\% = 0.97$$

$$\Pr[\text{Cancer}] = 1\% = 0.01$$

כמו כן, נסיק כי:

$$\Pr[\neg\text{Cancer}] = 1 - \Pr[\text{Cancer}] = 1 - 0.01 = 0.99$$

$$\Pr[1 | \neg\text{Cancer}] = 1 - \Pr[0 | \neg\text{Cancer}] = 1 - 0.97 = 0.03$$

$$\begin{aligned} \Pr[1] &= \Pr[1 | \text{Cancer}] \cdot \Pr[\text{Cancer}] + \Pr[1 | \neg\text{Cancer}] \cdot \Pr[\neg\text{Cancer}] = \\ &= 0.98 \cdot 0.01 + 0.03 \cdot 0.99 = 0.0395 \end{aligned}$$

לכן, ולפי כלל Bayes:

$$\begin{aligned} \Pr[\text{Cancer} | 1] &= \frac{\Pr[1 | \text{Cancer}] \cdot \Pr[\text{Cancer}]}{\Pr[1]} = \\ &= \frac{0.98 \cdot 0.01}{0.0395} = 0.248 = 24.8\% \end{aligned}$$

באופן מפתיע, אם הערכה מזהה מטופל כחולה בסרטן, בסיכוי של פחות מ-25% הוא אכן חולה, כלומר 3 מתוך 4 מטופלים שהערכה חוזה שהם חולים, הם למעשה בריאים. אם היינו רוצים להיות מדויקים יותר, יכולנו לבשר לכלל המטופלים שהם בריאים, ואז השגיאה הייתה 1% בלבד. הסיבה לכך שהשגיאה גדולה מאוד היא ההסתברות לסרטן בקרב כלל האוכלוסייה, שהיא קטנה מאוד - 1% בלבד.

**2.4 דוגמה: התפלגות נורמלית****2.4.1 התפלגות נורמלית**

**הגדרה** נאמר שמשתנה מקרי (מ"מ)  $Z$  מתפלג נורמלית  $N(\mu, \sigma^2)$ , ונסמן  $Z \sim N(\mu, \sigma^2)$  כאשר התוחלת שלו היא  $\mu$ , והשונות שלו היא  $\sigma^2$ .  
תכונות המ"מ  $Z \sim N(\mu, \sigma^2)$ :

$$E[Z] = \mu$$

$$\begin{aligned} \text{Var}[Z] &= E[(Z - E[Z])^2] = \\ &= E[Z^2] - E^2[Z] = \\ &= \sigma^2 \end{aligned}$$

$$\Pr[a \leq Z \leq b] = \int_a^b \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2} \cdot \left(\frac{\mu-x}{\sigma}\right)^2} dx$$

**2.4.2 תיאור הבעיה**

נניח כי  $\mu, \sigma \sim N(0, 1)$ , כלומר  $\mu$  ו- $\sigma$  הוגרלו מבעוד מועד על ידי מ"מ נורמלי סטנדרטי. נתון מ"מ  $Z \sim N(\mu, \sigma^2)$ . ננסה בעזרת למידה חישובית למצוא את  $\mu$  ו- $\sigma$ .  
נניח כי נתונות לנו  $n$  דגימות של המ"מ  $Z - z_1, z_2, \dots, z_n$ .  
לפי חוק Bayes:

$$\Pr[(\mu, \sigma) | z_1, z_2, \dots, z_n] = \frac{\Pr[z_1, z_2, \dots, z_n | (\mu, \sigma)] \cdot \Pr[(\mu, \sigma)]}{\Pr[z_1, z_2, \dots, z_n]}$$

כידוע:

$$\Pr[z_1, z_2, \dots, z_n | (\mu, \sigma)] = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2} \cdot \left(\frac{\mu-z_i}{\sigma}\right)^2}$$

$$\Pr[(\mu, \sigma)] = \frac{1}{\sqrt{2\pi}} \cdot e^{-\mu^2/2} \cdot \frac{1}{\sqrt{2\pi}} \cdot e^{-\sigma^2/2}$$

כמו כן, נשים לב כי מבחינתנו  $\Pr[z_1, z_2, \dots, z_n]$  הוא קבוע שרק מנרמל את הנוסחה, ולכן נשמיט אותו בחישובים שלנו.

**2.4.3 שיטת ML**

נאחזו מחפשים את ההשערה  $h_{ML}$ , שמוגדרת לפי:

$$h_{ML} = \arg \max_{h \in H} \Pr[\text{data} | h]$$

במקרה שלנו, אנחנו מחפשים את  $\mu$  ו- $\sigma$  שיביאו למקסימום את  $L(\mu, \sigma)$ , שמוגדר על ידי:

$$\begin{aligned} L(\mu, \sigma) &= \Pr[z_1, z_2, \dots, z_n | (\mu, \sigma)] = \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2} \cdot \left(\frac{\mu - z_i}{\sigma}\right)^2} \end{aligned}$$

נוציא log:

$$\log L(\mu, \sigma) = \sum_{i=1}^n -\frac{1}{2} \cdot \left(\frac{\mu - z_i}{\sigma}\right)^2 - \frac{n}{2} \log 2\pi - n \log \sigma$$

נגזור לפי  $\mu$ :

$$\begin{aligned} \frac{\partial}{\partial \mu} \log L(\mu, \sigma) &= \sum_{i=1}^n \frac{z_i - \mu}{\sigma} \cdot \frac{1}{\sigma} = 0 \\ \sum_{i=1}^n z_i &= n \cdot \mu \\ \mu &= \frac{1}{n} \cdot \sum_{i=1}^n z_i \end{aligned}$$

נגזור לפי  $\sigma$ :

$$\begin{aligned} \frac{\partial}{\partial \sigma} \log L(\mu, \sigma) &= \sum_{i=1}^n \frac{1}{\sigma} \cdot \left(\frac{z_i - \mu}{\sigma}\right)^2 - \frac{n}{\sigma} = 0 \\ \sum_{i=1}^n (z_i - \mu)^2 &= n \cdot \sigma^2 \\ \sigma^2 &= \frac{1}{n} \cdot \sum_{i=1}^n (z_i - \mu)^2 \end{aligned}$$

#### 2.4.4 שיטת MAP

אנחנו מחפשים את ההשערה  $h_{\text{MAP}}$ , שמוגרת לפי:

$$\begin{aligned} h_{\text{MAP}} &= \arg \max_{h \in H} \Pr[h | D] = \\ &= \arg \max_{h \in H} \frac{\Pr[D | h] \cdot \Pr[h]}{\Pr[D]} \end{aligned}$$

**הערה** מעתה,  $D = \text{data}$ .

במקרה שלנו, אנחנו מחפשים את  $\mu$  ו- $\sigma$  שיביאו למקסימום את  $L_{\text{MAP}}(\mu, \sigma)$ , שמוגדר על ידי:

$$L_{\text{MAP}}(\mu, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2} \cdot \left(\frac{z_i - \mu}{\sigma}\right)^2} \cdot \frac{1}{\sqrt{2\pi}} \cdot e^{-\mu^2/2} \cdot \frac{1}{\sqrt{2\pi}} \cdot e^{-\sigma^2/2}$$



נוציא log:

$$\log L_{\text{MAP}}(\mu, \sigma) = \sum_{i=1}^n -\frac{1}{2} \cdot \left(\frac{z_i - \mu}{\sigma}\right)^2 - \frac{n}{2} \log 2\pi - n \log \sigma - \frac{\mu^2}{2} - \frac{\sigma^2}{2} - \log 2\pi$$

נגזור לפי  $\mu$  ו- $\sigma$ :

$$\begin{aligned} \frac{\partial}{\partial \mu} \log L_{\text{MAP}}(\mu, \sigma) &= \sum_{i=1}^n \frac{z_i - \mu}{\sigma^2} - \mu = 0 \\ \frac{\partial}{\partial \sigma} \log L_{\text{MAP}}(\mu, \sigma) &= \sum_{i=1}^n \frac{(z_i - \mu)^2}{\sigma^3} - \frac{n}{\sigma} - \sigma = 0 \end{aligned}$$

קיבלנו מערכת משוואות על  $\mu$  ו- $\sigma$ . ננסה לפתור:

$$\begin{aligned} \frac{1}{n} \cdot \sum_{i=1}^n z_i &= \mu \cdot \left(1 + \frac{\sigma^2}{n}\right) \\ \frac{1}{n} \cdot \sum_{i=1}^n (z_i - \mu)^2 &= \sigma^2 \cdot \left(1 + \frac{\sigma^2}{n}\right) \end{aligned}$$

### 2.4.5 שיטת Posterior Bayes

נניח  $\mu \sim N(\eta, 1)$  ו- $Z \sim N(\mu, 1)$  (הנחנו כי  $\sigma = 1$ ). אז:

$$\begin{aligned} \Pr[\mu] &= \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2}(\mu - \eta)^2} \\ \Pr[z | \mu] &= \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2}(z - \mu)^2} \end{aligned}$$

ולפי חוק Bayes:

$$\begin{aligned} \Pr[\mu | z] &= \frac{\Pr[z | \mu] \cdot \Pr[\mu]}{\Pr[z]} = \\ &= \alpha \cdot e^{-\frac{1}{2}(\mu^2 - 2\mu\eta + \eta^2) - \frac{1}{2}(z^2 - 2\mu z + \mu^2)} = \\ &= \alpha \cdot \exp\left\{-\frac{1}{2} \cdot (2\mu^2 - 2\mu(\eta + z) + \eta^2 + z^2)\right\} = \\ &= \alpha \cdot \exp\left\{-\left(\mu - \frac{\eta + z}{2}\right)^2 + \underbrace{\left(\frac{\eta + z}{2}\right)^2 - \eta^2 - z^2}_{\text{Normalization}}\right\} \end{aligned}$$

כאשר  $\alpha$  הוא קבוע כלשהו (ניתן לחילוף מתוך החישוב לעיל).  
נגדיר כעת:

$$\begin{aligned} \hat{\mu} &= \frac{\eta + z}{2} \\ \hat{\sigma}^2 &= \frac{1}{2} \end{aligned}$$

---

<sup>2</sup>זו למעשה האמונה המוקדמת שלנו.

כעת, נוכל להניח כי  $Z \sim N(\hat{\mu}, \hat{\sigma}^2)$ , ונוכל להמשיך את התהליך. לחלופין, נניח שהגרלנו  $\mu \sim N(\eta, s^2)$  ו- $Z \sim N(\mu, \sigma^2)$ . נוכל להפעיל את התהליך על הנתונים  $z_1, z_2, \dots, z_n$  שהוגרלו לפי  $Z$ . נסמן:  $\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$ . כלל ההסקה שלנו יהיה:

$$\hat{\mu} = \frac{\frac{1}{s^2} \cdot \eta + \frac{n}{\sigma^2} \cdot \bar{z}}{\frac{1}{s^2} + \frac{n}{\sigma^2}}$$

$$\hat{\sigma}^2 = \left( \frac{1}{s^2} + \frac{n}{\sigma^2} \right)^{-1}$$

נשים לב שעבור ההשערה  $s = \sigma$ :

$$\hat{\mu} = \frac{\eta + n\bar{z}}{1 + n}$$

$$\hat{\sigma}^2 = \frac{\sigma^2}{1 + n}$$

מבחינה איכותית,  $\hat{\mu}$  מתקרב יותר ויותר לממוצע, ו- $\hat{\sigma}^2$  הולך ויורד שזה דבר טבעי ככל שיש יותר נתונים.

**הערה** ראינו כאן שתי אפשרויות: איך עושים את התהליך עבור נקודה אחת, ואיך עושים אותו עבור  $n$  נקודות. אולם, אם נעשה  $n$  פעמים את התהליך מנקודה לנקודה, נקבל בדיוק את אותה התוצאה שהיינו מקבלים מביצוע התהליך עבור  $n$  נקודות ישירות.

## Learning a Concept Class 2.5

נניח כי  $H$  מחלקת השערות, ו- $f \in H$  היא פונקציית המטרה. נניח כי  $\langle x, f(x) \rangle$  היא דוגמה.

נשים לב ש- $\Pr[h(x) = 1] \in \{0, 1\}$  עבור  $h \in H$  מסויימת, כי ידועה לנו מראש, ולכן אנחנו יכולים לבדוק אם  $h(x) = 1$  או  $h(x) = 0$ . נסמן:  $S = \{\langle x_i, b_i \rangle\}$ , כאשר  $b_i = f(x_i) \in \{0, 1\}$ . אזי:

$$\Pr[S | h] = \begin{cases} 0 & \exists i. b_i \neq h(x_i) \\ \Pr[x_1, x_2, \dots, x_n] & \forall i. b_i = h(x_i) \end{cases}$$

**הגדרה** נאמר ש- $h \in H$  עקבית אם היא מסווגת את כל הנקודות ב- $S$  בצורה נכונה, כלומר  $\forall i. b_i = h(x_i)$ .

נגדיר את  $H' \subseteq H$  לפי  $H' = \{h \in H \mid h \text{ is consistent}\}$ . נסתכל על כלל הסקה שיש לנו:

ML יחזיר השערה עקבית כלשהי.

MAP יחזיר השערה עקבית עם המשקל הפוסטריורי הגבוה ביותר, כלומר את ההשערה העקבית שהכי האמנו בה.

Bayes support יהיה  $H'$  ומשקול מחדש<sup>3</sup>. המשקול המחדש יבוצע לפי הנוסחה:

$$B(y) = \sum_{h \in H'} \frac{h(y) \cdot \Pr[h]}{\Pr[H']}$$

## 2.6 דוגמה: Biased Coins

נניח כי מטילים מטבע  $m$  פעמים, ורואים  $k$  הצלחות. רוצים להעריך את ההסתברות  $p$  של המטבע. נחשב את  $p$  בעזרת ML:

$$\Pr[(k, m) | p] = \binom{m}{k} \cdot p^k \cdot (1-p)^{m-k}$$

לאחר חישוב מקבלים כי  $p = k/m$ . איך עושים את זה? מוציאים  $\log$  ומחלקים ב- $m$ , ואז אפשר לגזור. לא נראה כאן את החישוב.

### 2.6.1 חוק Laplace

נניח התפלגות prior אחידה (Uniform). כלומר, התפלגות לכל המטבעות האפריות היא אחידה:

$$\Pr[p \leq \theta] = \int_0^\theta dp = \theta$$

נחשב את ההסתברות ל- $k$  הצלחות מ- $m$  הטלות:

$$\begin{aligned} \Pr[(k, m)] &= \int_0^1 \Pr[k | p] \cdot \Pr[p] dp = \\ &= \int_0^1 \binom{m}{k} x^k (1-x)^{m-k} dx = \\ &= \left[ \binom{m}{k} \cdot \frac{x^{k+1}}{k+1} \cdot (1-x)^{m-k} \right]_0^1 \\ &\quad + \int_0^1 \binom{m}{k} \cdot \frac{x^{k+1}}{k+1} \cdot (1-x)^{m-k-1} \cdot (m-k) dx = \\ &= \int_0^1 \binom{m}{k+1} \cdot x^{k+1} \cdot (1-x)^{m-k+1} dx = \\ &= \int_0^1 \Pr[k+1 | p] \cdot \Pr[p] dp = \Pr[(k+1, m)] \end{aligned}$$

נשים לב שהשתמשנו כאן במעברים בזהות הבאה:

$$\binom{n}{k} \cdot \frac{n-k}{k+1} = \binom{n}{k+1}$$

<sup>3</sup>למשל, אם  $H = \{h_1, h_2, h_3, h_4\}$  וכך:  $\Pr[h_i] = i/10$ .  $\forall 1 \leq i \leq 4$ . נניח כי  $H' = \{h_2, h_3\}$ . נצטרך למשקל מחדש את  $h_2$  ו- $h_3$ , כך ש- $\Pr[h_2] + \Pr[h_3] = 1$ .

נשים לב:

$$\begin{aligned} \Pr [(k, m)] &= \int_0^1 p^k \cdot (1-p)^{m-k} dp = \\ &= \frac{1}{\binom{m}{k}} \cdot \frac{1}{m+1} \end{aligned}$$

נחשב:

$$\begin{aligned} E [p | (k, m)] &= \int_0^1 p \cdot \frac{\Pr [(k, n) | p] \cdot \Pr [p]}{\Pr [(k, m)]} dp = \\ &= \int_0^1 p \cdot \frac{p^k \cdot (1-p)^{m-k}}{\frac{1}{m+1} \cdot \frac{1}{\binom{m}{k}}} dp = \\ &= (m+1) \cdot \binom{m}{k} \cdot \int_0^1 p^{k+1} \cdot (1-p)^{m-k} dp = \\ &= (m+1) \cdot \binom{m}{k} \cdot \frac{1}{m+2} \cdot \frac{1}{\binom{m+1}{k+1}} = \\ &= \frac{m+1}{m+2} \cdot \frac{k+1}{m+1} = \\ &= \frac{k+1}{m+2} \end{aligned}$$

למעשה, התיקון של Laplace מוסיף לנו עוד שתי דגימות, כאשר רק באחת מהן יש הצלחה. למה זה יותר טוב? אם יש הרבה הטלות, זה עדיין קרוב למה שהיינו חושבים. אולם אם יש מעט הטלות, אז עדיין יש מקום למקרי הקיצון.

### 2.6.2 פונקציות Loss

למעשה, לא נוכל אף פעם לחזות מה יקרה בהטלת המטבע הבאה. לכן, נכנס לכאן העניין של ההפסד (Loss).

כאשר מדברים על פונקציות Loss, יש לקחת בחשבון שני מרכיבים להפסד:

1. Bayes Risk - זה ההפסד הבלתי-נמנע, גם אם אנחנו מכירים את כל הפרטים במערכת. למשל, בבעיית המטבע, לא נוכל אף פעם לחזות במדויק מה תהיה תוצאת ההטלה הבאה, אלא רק בהסתברות מסויימת.

2. Regret - ההפסד בגלל הערכה לא נכונה של המערכת.

**תזכורת** פונקציית log loss מוגדרת באופן הבא:

$$\ell_{\log}(x) = \begin{cases} \log \frac{1}{p} & f(x) = 1 \\ \log \frac{1}{1-p} & f(x) = 0 \end{cases}$$

עבור  $\theta$ -bias, ה־log loss הצפוי הוא:

$$\ell_{\log} = \theta \cdot \log \frac{1}{p} + (1 - \theta) \cdot \log \frac{1}{1 - p}$$

כמובן שהערך המינימלי מתקבל עבור  $p = \theta$ . ואם באמת  $p = \theta$ , אז נסמן:

$$H(\theta) = \theta \cdot \log \frac{1}{\theta} + (1 - \theta) \cdot \log \frac{1}{1 - \theta}$$

לערך זה נקרא האנטרופיה של  $\theta$ . זהו ההפסד שנובע מהמערכת עצמה, ולא מכך שאיננו מכירים את המערכת, ולכן זהו ה־Bayes Risk. נבדוק כמה אנחנו רחוקים מה־Bayes Risk (כלומר מהו ההפסד הנוסף, ה־Regret) כשאנחנו משתמשים בחוק Laplace: נניח כי יש לנו  $T$  דגימות בזמן. אז:

$$\begin{aligned} E[\log \text{loss}] &= \\ &= \int_0^1 \sum_{m=1}^T \sum_{k=0}^m \left[ \theta \cdot \log \frac{m+2}{k+1} + (1-\theta) \cdot \log \frac{m+2}{m-k+1} \right] \\ &\quad \binom{m}{k} \cdot \theta^k \cdot (1-\theta)^{m-k} d\theta = \\ &= \sum_{m=1}^T \sum_{k=0}^m \binom{m}{k} \cdot \log \frac{m+2}{k+1} \cdot \int_0^1 \theta^{k+1} \cdot (1-\theta)^{m-k} d\theta + \\ &\quad \sum_{m=1}^T \sum_{k=0}^m \binom{m}{k} \cdot \log \frac{m+2}{m-k+1} \cdot \int_0^1 \theta^k \cdot (1-\theta)^{m-k+1} d\theta = \\ &= \sum_{m=1}^T \sum_{k=0}^m \left[ \frac{1}{m+1} \cdot \frac{k+1}{m+2} \cdot \log \frac{m+2}{k+1} + \frac{1}{m+1} \cdot \frac{m-k+1}{m+2} \cdot \log \frac{m+2}{m-k+1} \right] = \\ &= \sum_{m=1}^T \sum_{k=0}^m \frac{1}{m+1} \cdot H\left(\frac{k+1}{m+1}\right) \sim \\ &\sim T \cdot \int_0^1 H(\theta) d\theta + O\left(\underbrace{\sum_{m=1}^T \frac{1}{m}}_{\log T}\right) \end{aligned}$$

ננסה לחסום את  $T \cdot \int_0^1 H(\theta) d\theta$ :

$$\sum_{i=1}^{m/2} \frac{1}{m} \cdot H\left(\frac{i-1}{m}\right) \leq \int_0^1 H(\theta) d\theta \leq \sum_{i=1}^{m/2} \frac{1}{m} \cdot H\left(\frac{i}{m}\right)$$

$$\begin{aligned} \sum_{i=1}^{m/2} \frac{1}{m} \cdot H\left(\frac{i}{m}\right) - \sum_{i=1}^{m/2} \frac{1}{m} \cdot H\left(\frac{i-1}{m}\right) &= \frac{1}{m} \cdot \sum_{i=1}^{m/2} \left[ H\left(\frac{i}{m}\right) - H\left(\frac{i-1}{m}\right) \right] = \\ &= \frac{1}{m} \cdot (1 - 0) = \frac{1}{m} \xrightarrow{m \rightarrow \infty} 0 \end{aligned}$$

$x_1$	$x_2$	...	$x_n$	$C$
0	1		1	+1
1	0		0	-1
⋮	⋮		⋮	⋮
0	1		0	+1

טבלה 2.1: דוגמה למדגם עבור סיווג בייסיאני למרחב בינארי

כלומר, ככל ש- $m$  גדל, כך  $\int_0^1 H(\theta) d\theta \cdot T$  הולך ומתקרב ל- $H(\frac{i}{m}) \cdot T \sum_{i=1}^{m/2} \frac{1}{m}$ , שהוא ה-Bayes Risk. לכן, הראנו שבאמצעות הפעלת חוק Laplace, קיבלנו את ה-loss האופטימלי (ה-Bayes Risk), עם תוספת לוגריתמית במספר הטלות המטבע (T).

## 2.7 Naïve Bayes

### 2.7.1 סיווג בייסיאני: מרחב בינארי

נתון המצב הבא: יש לנו שתי מחלקות +1 ו-1, וכל דוגמה מתוארת על ידי  $n$  מאפיינים, כאשר  $x_i$  ( $1 \leq i \leq n$ ) הוא משתנה בינארי, שערכו 0 או 1. דוגמה לקלט כזה ניתן למצוא בטבלה 2.1.

רוצים לבנות השערה  $h: \{0, 1\}^n \rightarrow \{+1, -1\}$ . לפי כלל Bayes:

$$\Pr[C = +1 | x_1, x_2, \dots, x_n] = \frac{\Pr[x_1, x_2, \dots, x_n | C = +1] \cdot \Pr[C = +1]}{\Pr[x_1, x_2, \dots, x_n]}$$

קל להעריך את  $\Pr[C = +1]$  מהנתונים. אבל איך נעריך את  $\Pr[x_1, x_2, \dots, x_n | C = +1]$ ? Naïve Bayes מבוסס על הנחת אי-תלות:

$$\Pr[x_1, x_2, \dots, x_n | C = +1] = \prod_{i=1}^n \Pr[x_i | C = +1]$$

כל מאפיין  $x_i$  הוא ב"ת באחרים ברגע שאנחנו יודעים את הערך של  $C$ . לכן, לכל  $1 \leq i \leq n$ , יש לנו שני מאפיינים:

$$\begin{aligned} \theta_{i,+1} &= \Pr[x_i = 1 | C = +1] \\ \theta_{i,-1} &= \Pr[x_i = 1 | C = -1] \end{aligned}$$

קיבלנו כאן  $2n$  פרמטרים בלתי-תלויים (ב"ת).

**תזכורת**

אי שוויון Markov: אם  $X \geq 0$  מ"מ, אזי  $\Pr[X \geq \lambda] \leq \frac{E[X]}{\lambda}$ .

אי שוויון Cheviechev: אם  $X \geq 0$  מ"מ, אזי  $\Pr[X \geq \lambda] \leq \frac{E[X^2]}{\lambda^2}$ .

אי שוויון Chernoff: אם  $X_1, X_2, \dots, X_m$  מ"מ ב"ת כך ש- $E[X_i] = p$ , אזי  $\Pr\left[\left|\frac{1}{m} \cdot \sum_{i=1}^m X_i - p\right| \geq \lambda\right] \leq e^{-2\lambda^2 m}$ .

## 2.7.2 פענוח של Naïve Bayes

לפי Bayes ו-MAP, אנחנו צריכים להעריך את  $\Pr[C = +1 | x_1, x_2, \dots, x_n]$  לעומת  $\Pr[C = -1 | x_1, x_2, \dots, x_n]$ . קל לעשות זאת בעזרת  $\log$  וחלוקה (והשוואה ל-0):

$$\begin{aligned} & \log \frac{\Pr[C = +1 | x_1, x_2, \dots, x_n]}{\Pr[C = -1 | x_1, x_2, \dots, x_n]} = \\ &= \log \frac{\Pr[x_1, x_2, \dots, x_n | C = +1] \cdot \Pr[C = +1]}{\Pr[x_1, x_2, \dots, x_n | C = -1] \cdot \Pr[C = -1]} = \\ &= \log \frac{\Pr[C = +1]}{\Pr[C = -1]} + \log \prod_{i=1}^n \frac{\Pr[x_i | C = +1]}{\Pr[x_i | C = -1]} = \\ &= \log \frac{\Pr[C = +1]}{\Pr[C = -1]} + \sum_{i=1}^n \log \frac{\Pr[x_i | C = +1]}{\Pr[x_i | C = -1]} \end{aligned}$$

לכן, הסקנו כי

$$\log \frac{\Pr[C = +1 | x_1, x_2, \dots, x_n]}{\Pr[C = -1 | x_1, x_2, \dots, x_n]} = \log \frac{\Pr[C = +1]}{\Pr[C = -1]} + \sum_{i=1}^n \log \frac{\Pr[x_i | C = +1]}{\Pr[x_i | C = -1]}$$

כלומר, כל משפיע על החיזוי:

- אם  $\Pr[x_i | C = +1] = \Pr[x_i | C = -1]$ , אז ל- $x_i$  אין השפעה על החיזוי.
- אם  $\Pr[x_i | C = -1] = 0$ , אז  $x_i$  משפיע על שאר הקולות (דומה להטלת וטו).
- באופן דומה, אם  $\Pr[x_i | C = +1] = 0$ .

נסמן:

$$\begin{aligned} w_i &= \log \frac{\Pr[x_i = 1 | C = +1]}{\Pr[x_i = 1 | C = -1]} - \log \frac{\Pr[x_i = 0 | C = +1]}{\Pr[x_i = 0 | C = -1]} \\ b &= \log \frac{\Pr[C = +1]}{\Pr[C = -1]} + \sum_{i=1}^n \log \frac{\Pr[x_i = 0 | C = +1]}{\Pr[x_i = 0 | C = -1]} \end{aligned}$$

כלל ההחלטה שלנו יהיה  $\text{sign}(b + \sum_{i=1}^n x_i w_i)$ .

## 2.7.3 התפלגות נורמלית

השלב הבא ב-Naïve Bayes אומר ש- $\Pr[x_i | C] \sim N(\mu_{i,C}, \sigma_i)$ .<sup>4</sup> נעשה את אותו החישוב:

$$\begin{aligned} & \log \frac{\Pr[C = +1 | x_1, x_2, \dots, x_n]}{\Pr[C = -1 | x_1, x_2, \dots, x_n]} = \\ &= \log \frac{\Pr[C = +1]}{\Pr[C = -1]} + \sum_{i=1}^n \log \frac{\Pr[x_i | C = +1]}{\Pr[x_i | C = -1]} \end{aligned}$$

<sup>4</sup>חשוב לשים לב ש- $\mu_{i,C}$  תלוי ב- $C$ , ואילו  $\sigma_i$  אינו תלוי ב- $C$ .

$$\begin{aligned}
\log \frac{\Pr[x_i | C = +1]}{\Pr[x_i | C = -1]} &= \log \frac{e^{-\frac{1}{2} \cdot \left(\frac{\mu_{i,+1} - x_i}{\sigma_i}\right)^2}}{e^{-\frac{1}{2} \cdot \left(\frac{\mu_{i,-1} - x_i}{\sigma_i}\right)^2}} = \\
&= -\frac{1}{2} \cdot \left(\frac{\mu_{i,+1} - x_i}{\sigma_i}\right)^2 + \frac{1}{2} \cdot \left(\frac{\mu_{i,-1} - x_i}{\sigma_i}\right)^2 = \\
&= \frac{1}{2 \cdot \sigma^2} \cdot (\mu_{i,+1} + \mu_{i,-1} - 2x_i) \cdot (\mu_{i,-1} - \mu_{i,+1}) = \\
&= \frac{1}{2} \cdot \frac{\mu_{i,-1} - \mu_{i,+1}}{\sigma_i} \cdot \frac{\mu_{i,+1} + \mu_{i,-1}}{2} - x_i
\end{aligned}$$



## פרק 3

# מודל ה-PAC<sup>1</sup>

### 3.1 מודל ה-PAC

ראשי התיבות PAC פירושו *Probably Approximately Correct*. המטרה של המודל היא למצוא השערה  $h$ , שבהסתברות גבוהה (Probably) היא מדוייקת (Approximately Correct). מבנה השיעור:

1. דוגמה ללימוד PAC.
2. מודל פורמלי ומחלקת השערות סופית.
3. Occam Razor ודוגמאות.

### 3.2 דוגמה אינטואיטיבית

נניח כי בן-אדם טיפוסי מקיים את החיתוך של שתי התכונות הבאות:

$$\begin{aligned} 60 &\leq \text{Weight} \leq 90 \\ 1.60 &\leq \text{Height} \leq 1.90 \end{aligned}$$

אזי  $R$  היא מחלקת בני האדם. נשים לב לכך ש- $R$  הוא מלבן על המישור. אנחנו נמצא מחלקה  $R'$  שהיא קירוב של  $R$ . במודל PAC אין הנחה על התפלגות הדוגמאות. ההנחה היחידה שלנו היא שקיימת התפלגות, והיא נדגמת באופן i.i.d. נתאר את הבעיה:

- קלט: אוסף  $S$  של דוגמאות מסווגות.
- פלט:  $R'$  (מלבן).
- מטרה: משקל קטן עבור  $D(R \Delta R')$  (חיסור סימטרי).

---

<sup>1</sup>שיעור שהתקיים בתאריך 04.11.2012.

**3.2.1 מציאת השערה טובה**

על השגיאה ניתן להסתכל באופן הבא:

$$R \Delta R' = \underbrace{(R \setminus R')}_{\text{False-Positive}} \cup \underbrace{(R' \setminus R)}_{\text{False-Negative}}$$

המטרה שלנו היא למצוא השערה  $R'$ , כך שבהסתברות  $1 - \delta$ :

$$\Pr[\text{error}] = D(R \Delta R') \leq \varepsilon$$

כאשר  $R$  היא פונקציית המטרה.

**3.2.2 אופן הלימוד**

אם נסתכל על המישור, יש שתי בחירות אינטואיטיביות:

- $R_{\min}$  - המלבן החוסם הקטן ביותר.
- $R_{\max}$  - המלבן החוסם הגדול ביותר.

איך נחשב את  $R_{\min}$  (לדוגמה) בהינתן קבוצת דגימות  $S = \{\langle (x_1, y_1), b_1 \rangle, \dots, \langle (x_m, y_m), b_m \rangle\}$  נסמן:

$$m_x = x_1 = M_x$$

$$m_y = y_1 = M_y$$

נעבור על  $\langle (x_i, y_i), b_i \rangle$ . אם  $b_i = +$ :

- אם  $m_x > x_i$ , אז נגדיר  $m_x = x_i$ .
- אם  $M_x < x_i$ , אז נגדיר  $M_x = x_i$ .
- אם  $m_y > y_i$ , אז נגדיר  $m_y = y_i$ .
- אם  $M_y < y_i$ , אז נגדיר  $M_y = y_i$ .

כעת,  $R_{\min} = \{(x, y) \mid m_x \leq x \leq M_x \wedge m_y \leq y \leq M_y\}$ . נשים לב שזהו אלגוריתם יעיל (פולינומיאלי) ב- $m$ , מספר הדגימות שלנו.

**3.2.3 מספר הדגימות**

ננתח את הבעיה באופן פורמלי. נתונים לנו:

- $\varepsilon$  - דיוק.
- $\delta$  - ביטחון.

אנחנו מחפשים אלגוריתם  $A$ , כך שעבור מדגם בגודל  $m(\varepsilon, \delta)$ , מחזיר מלבן  $R'$ , כך שבהסתברות  $1 - \delta$ :

$$\Pr[\text{error}] \leq \varepsilon$$

ננתח את  $m(\varepsilon, \delta)$  עבור האלגוריתם  $A$  שמחזיר את  $R_{\min}$ . נשים לב ש- $R_{\min} \subseteq R$ , ולכן:

$$\Pr[\text{error}] = D(R \Delta R_{\min}) = D(R \setminus R_{\min})$$

לכן, נשים לב שהשגיאה האפשרית שלנו נמצאת כולה ב- $R$ . נחלק את השגיאה שלנו לארבעה מלבנים:  $T'_1$  (החלק שמעל  $R_{\min}$ ),  $T'_2$  (החלק שמשמאל ל- $R_{\min}$ ),  $T'_3$  (החלק שמתחת ל- $R_{\min}$ ) ו- $T'_4$  (החלק שמימין ל- $R_{\min}$ ). אזי:

$$R \Delta R_{\min} = \bigcup_{i=1}^4 T'_i$$

$$D(R \Delta R_{\min}) \leq \sum_{i=1}^4 D(T'_i)$$

לכן, אם נראה כי  $\forall i. D(T'_i) \leq \varepsilon/4$ , אז סיימנו, כי:

$$D(R \Delta R_{\min}) \leq 4 \cdot \varepsilon/4 = \varepsilon$$

הבעיה הקונסטרוקטיבית שלנו במצב הזה היא שהגדרנו את  $T'_i$  רק אחרי שכבר ראינו את המדגם ובנינו את  $R_{\min}$ . החוכמה היא להגדיר מאורעות שאינם תלויים במדגם. רק אז נוכל לחשב מה המשקל שלהם.

נשים לב לכך שמכיוון שבהכרח  $R_{\min} \subseteq R$ , כל השגיאה תהיה תמיד ב- $R \setminus R_{\min}$ . לכן, נחפש  $T_1, T_2, T_3, T_4$  באופן הבא:

- $T_1$  הוא תת-מלבן של  $R$  הצמוד לדופן העליונה של  $R$ .
- $T_2$  הוא תת-מלבן של  $R$  הצמוד לדופן השמאלית של  $R$ .
- $T_3$  הוא תת-מלבן של  $R$  הצמוד לדופן התחתונה של  $R$ .
- $T_4$  הוא תת-מלבן של  $R$  הצמוד לדופן הימנית של  $R$ .
- לכל  $1 \leq i \leq 4$  מתקיים  $D(T_i) = \varepsilon/4$ .

נשים לב לכך שהבחירה של  $\{T_i\}_{i=1}^4$  אינה תלוייה במדגם. נרצה לחשב את ההסתברות ש- $T'_i \subseteq T_i$ , כלומר, שבהסתברות  $1 - \delta$ ,  $\forall i. T'_i \subseteq T_i$ . נסתכל על  $T_1$ : מתי  $T'_1 \subseteq T_1$ ? אם קיימת ב- $S$  דוגמה  $\langle (x, y), + \rangle$  כך ש- $(x, y) \in T_1$ , אזי  $T'_1 \subseteq T_1$  לפי הבנייה:

$$\Pr[(x, y) \notin T_1] = 1 - \frac{\varepsilon}{4}$$

מכיוון שהדגימות של  $S$  הן i.i.d, נסיק כי:

$$\Pr[\forall (x_i, y_i) \in S. (x_i, y_i) \notin T_1] = \left(1 - \frac{\varepsilon}{4}\right)^m$$

את אותם החישובים ניתן לעשות גם עבור  $\{T_i\}_{i=2}^4$ . לכן:

$$\begin{aligned} \Pr[\text{error}] &= \Pr[\exists i. \forall (x, y) \in S. (x, y) \notin T_i] \\ &\leq \sum_{i=1}^4 \Pr[\forall (x, y) \in S. (x, y) \notin T_i] \\ &\leq 4 \cdot \left(1 - \frac{\varepsilon}{4}\right)^m \leq 4e^{-\frac{\varepsilon}{4}m} < \delta \end{aligned}$$

המעבר האחרון באי-השוויון בוצע לפי הזהות  $1 - x \leq e^{-x}$ . לסיום, נסיק כי:

$$m > \frac{4}{\varepsilon} \cdot \ln \frac{4}{\delta}$$

### 3.3 הצגה פורמלית של מודל ה-PAC

#### 3.3.1 הקדמה

- המטרה שלנו היא ללמוד השערה מתוך קבוצה ידועה מראש של השערות.
- הסביבה סטוכסטית.
- הדגימות נדגמות מהתפלגות i.i.d.
- ההתפלגות על הלמידה (*train*) זהה להתפלגות על הבדיקה (*test*).
- הפתרון צריך להיות יעיל: גודל המדגם ביחס לזמן החישוב צריך להיות פולינומיאלי ב- $\frac{1}{\varepsilon}$  ו- $\frac{1}{\delta}$ .

#### 3.3.2 הגדרת מודל ה-PAC

יהי  $X$  מרחב הדוגמאות, ותהי  $D$  התפלגות מעל  $X$ . נאמר כי  $C \subseteq \{c \mid c : X \rightarrow \{0, 1\}\}$  היא פחלקת פונקציות המטרה. תהי  $c_t \in C$  פונקציית המטרה שלנו. תהי  $H \subseteq C$  מחלקת ההשערות שלנו. תהי  $h \in H$  נגדיר את השגיאה של ההשערה  $h$ :

$$\text{error}(h) = \Pr_D[h(x) \neq c_t(x)] = D(h \Delta c_t)$$

כמו כן, נגדיר את  $EX(c_t, D)$  להיות אורקל (*Oracle*), המחזיר דגימה  $x \in X$  שנדגמה לפי  $D$  וסיווג  $c_t(x)$ :  $\langle x, c_t(x) \rangle$ .

**הגדרה** נאמר כי  $C$  נלפדת PAC על ידי  $H$  אם קיים אלגוריתם  $A$  כך שלכל פונקציית מטרה  $c_t \in C$ , לכל התפלגות  $D$  מעל המרחב  $X$ , ולכל פרמטרים  $\varepsilon, \delta > 0$ , אם נותנים ל- $A$  גישה ל- $EX(c_t, D)$ , אזי בהסתברות  $1 - \delta$ , יחזיר השערה  $h \in H$  כך ש:

- אם  $c_t \in H$  (*realizable*), אזי:

$$\text{error}(h) \leq \varepsilon$$

• אם  $c_t \notin H$  (non-realizable), אזי:

$$\text{error}(h) \leq \varepsilon + \min_{\tilde{h} \in H} \text{error}(\tilde{h})$$

נאמר כי  $C$  נלפזת זיעילות אם  $A$  רץ בזמן פולינומיאלי ב- $\frac{1}{\varepsilon}$ ,  $\ln \frac{1}{\delta}$ , וגם ב- $n$  ו- $\ell$ , כאשר  $n$  הוא גודל הדוגמה (בביטים) ו- $\ell$  הוא גודל פונקציית המטרה  $c_t$  (גם כן בביטים).

### 3.4 מחלקות השערות סופיות

נדון במקרה בו  $H$  סופית.

#### 3.4.1 המקרה $c_t \in H$

**הגדרה** נאמר כי השערה  $h \in H$  עקבית אם  $\forall x \in S. h(x) = c_t(x)$ . במקרה שלנו, מובטח לנו שקיימת  $h$  עקבית. נניח כי קיימת  $h \in H$  עקבית כך ש- $\text{error}(h) > \varepsilon$ . אזי נקראת מאורע רע.

$$\Pr[\forall i. h(x_i) = c_t(x_i) \mid \text{error}(h) > \varepsilon] \leq (1 - \varepsilon)^m \leq e^{-\varepsilon m}$$

$$\Pr[A \text{ returns } h \wedge \text{error}(h) > \varepsilon] \leq \delta$$

$$\Pr[\exists h \in H. \text{error}(h) > \varepsilon \wedge \forall i. h(x_i) = c_t(x_i)] \leq |H| \cdot e^{-\varepsilon m} \leq \delta$$

נחלץ את  $m$ :

$$\begin{aligned} |H| \cdot e^{-\varepsilon m} &\leq \delta \\ e^{-\varepsilon m} &\leq \frac{\delta}{|H|} \\ -\varepsilon m &\leq \ln \frac{\delta}{|H|} \\ m &\geq \frac{1}{\varepsilon} \ln \frac{|H|}{\delta} \end{aligned}$$

#### 3.4.2 המקרה $c_t \notin H$

נסמן:

$$h^* = \arg \min_{h \in H} \text{error}(h)$$

אז, לכל  $h \in H$ :

$$\text{error}(h) \geq \text{error}(h^*) > 0$$

נסמן:

$$\beta = \text{error}(h^*)$$

המטרה שלנו תהיה:

$$\text{error}(h) \leq \beta + \varepsilon$$

נגדיר:

$$\widehat{\text{error}}(h) = \frac{1}{m} \sum_{i=1}^m I(h(x_i) \neq c_i(x_i))$$

הבחירה הטבעית שלנו תהיה:

$$\bar{h} = \arg \min_{h \in H} \widehat{\text{error}}(h)$$

בחירה זו נקראת גם ERM<sup>2</sup>.

נבחר מדגם מספיק גדול בגודל  $m$  כך שבהסתברות  $1 - \delta$  מתקיים:

$$\forall h \in H. |\widehat{\text{error}}(h) - \text{error}(h)| \leq \frac{\varepsilon}{2}$$

כעת, נקבל כי:

$$\begin{aligned} \text{error}(\bar{h}) &\leq \widehat{\text{error}}(\bar{h}) + \frac{\varepsilon}{2} \leq \widehat{\text{error}}(h^*) + \frac{\varepsilon}{2} \\ &\leq \text{error}(h^*) + \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \beta + \varepsilon \end{aligned}$$

לסיים, בהינתן  $h \in H$ :

$$\begin{aligned} \Pr [|\widehat{\text{error}}(h) - \text{error}(h)| \geq \varepsilon/2] &\leq 2e^{-2(\varepsilon/2)^2 m} \cdot |H| \leq \delta \\ e^{-2(\varepsilon/2)^2 m} &\leq \frac{\delta}{2|H|} \\ -2\left(\frac{\varepsilon}{2}\right)^2 m &\leq \ln \frac{\delta}{2|H|} \\ m &\geq \frac{8}{\varepsilon^2} \cdot \ln \frac{2 \cdot |H|}{\delta} \end{aligned}$$

**הערה** למה במקרה הראשון (realizable) מתנהג ביחס הפוך ל- $\varepsilon$ , ואילו במקרה השני (non-realizable) מתנהג ביחס הפוך ל- $\varepsilon^2$ ? ההבדל הוא שבמקרה הראשון, דגימה אחת מספיקה כדי לפסול השערה. לעומת זאת, במקרה השני, גם אם  $h$  אינה עקבית, היא עלולה להיות ההשערה הטובה ביותר.

<sup>2</sup>ראשי תיבות של Empirical Risk Management.

**אלגוריתם 1.3** ELIM ללמידת Boolean Disjunctions

בהינתן דגימה בגודל  $m$ , מתחילים עם  $L$ .  
 לכל דגימה שלילית  $\langle x, - \rangle$ , מוציאים את כל הליטרלים החיוביים.  
 כעת, כל הדוגמאות השליליות מסווגות נכון. גם כל הדוגמאות החיוביות מסווגות נכון, כי  $c_t \subseteq L_{\text{final}}$ .

**3.4.3 דוגמה - למידת Boolean Disjunctions**

נתונים  $n$  משתנים בולאנים  $x_1, \dots, x_n$ , וקבוצת ליטרלים:

$$L = \{x_1, \bar{x}_1, \dots, x_n, \bar{x}_n\}$$

נגדיר *disjunction* על ידי ביטוי OR של ליטרלים. למשל:  $x_1 \vee x_3 \vee \bar{x}_5$ .  
 נגדיר את  $C = H$  להיות מחלקת ה-disjunctions. אזי:  $|C| = |H| = 3^n$  (כי כל משתנה בולאני יכול להופיע בשני ליטרלים או לא להופיע, ששקול להופעת שני הליטרלים המתאימים).  $X = \{0, 1\}^n$  ופונקציית המטרה שלנו היא:

$$c_t = \bigvee_{\substack{j \in S \\ \ell_j \in L}} \ell_j$$

נשים לב לעובדה שבהינתן דגימה שלילית  $\langle x, - \rangle$ , ניתן לפסול את כל הליטרלים החיוביים שלה. למשל,  $\langle 001, - \rangle$  אומר לנו כי  $\bar{x}_1, \bar{x}_2, x_3 \notin c_t$ . זה הרעיון שעומד בבסיסו של האלגוריתם ELIM (אלגוריתם 1.3).  
 לפי חישוב קודם:

$$m \geq \frac{1}{\varepsilon} \ln \frac{|H|}{\delta} = \frac{n \ln 3}{\varepsilon} + \frac{1}{\varepsilon} \ln \frac{1}{\delta}$$

**3.4.4 דוגמה - למידת Parity**

במקרה הזה:

$$X = \{0, 1\}^n$$

$$c_t = \bigoplus_{j \in S} x_j$$

למשל,  $c_t = x_1 \oplus x_3 \oplus x_5$ . אזי  $|C| = 2^n$ .  
 אפשר להסתכל על הבעיה כאל בעיית פתרון מערכת משוואות. נגדיר את  $z_i$  כך ש-  
 $z_i = 1$  פירושו  $i \in S$ . אזי נוכל להסיק מערכת משוואות בצורה פשוטה. למשל, מ- $\langle 0110, + \rangle$   
 נסיק את המשוואה  $z_2 + z_3 = 1$ , ומ- $\langle 1111, - \rangle$  נסיק את המשוואה  $z_1 + z_2 + z_3 + z_4 = 0$   
 (כמובן שהמשוואות הן מעל  $\mathbb{Z}_2$ ).

קיבלנו מערכת משוואות אותה צריך לפתור. נסיק כי:

$$m \geq \frac{1}{\varepsilon} \ln \frac{|H|}{\delta} = \frac{1}{\varepsilon} \ln \frac{2^n}{\delta}$$

$$= \frac{n \ln 2}{\varepsilon} + \frac{1}{\varepsilon} \ln \frac{1}{\delta}$$

## Occam Razor 3.5

Entities should not be multiplied unnecessarily

(מתוך דבריו של William Occam בשנת 1320).

המשמעות שלנו למשפט שלו היא שניתן ל- $H$  לגדול יחד עם גודל המדגם  $m$ .

**הגדרה**  $(\alpha, \beta)$  הוא אלגוריתם Occam ללימוד מחלקה  $C$  על ידי מחלקה  $H$  אם  $\alpha \geq 0$  ו- $0 \leq \beta < 1$ , ובהינתן דגימה בגודל  $m$ , האלגוריתם מוציא השערה  $h \in H$  כך ש- $h$  עקבית, וגודל  $h$  חסור על ידי  $n^\alpha m^\beta$  (כאשר  $n$  הוא גודל דוגמה בודדת בביטים).

## 3.5.1 אלגוריתמי Occam ומוזל ה-PAC

**משפט** יהי  $A$  אלגוריתם Occam עבור  $C$  המשתמש ב- $H$ . אזי  $A$  הוא אלגוריתם PAC, כאשר:

$$m \geq \left( \frac{n^\alpha}{\epsilon} \ln 2 \right)^{\frac{1}{1-\beta}} + \frac{2}{\epsilon} \ln \frac{1}{\delta}$$

**הוכחה** נקבע את  $m$  ו- $n$ . אזי  $A$  מחזיר השערה  $h$  עם לכל היותר  $n^\alpha m^\beta$  ביטים. לכן, מספר ההשערות האפשריות חסום על ידי  $2^{n^\alpha m^\beta}$ .  $|H| \leq 2^{n^\alpha m^\beta}$ . לכן:

$$\begin{aligned} m &\geq \frac{1}{\epsilon} \ln \frac{|H|}{\delta} \geq \frac{n^\alpha m^\beta \ln 2}{\epsilon} + \frac{1}{\epsilon} \ln \frac{1}{\delta} \\ m &\geq \max \left\{ \frac{2n^\alpha m^\beta \ln 2}{\epsilon}, \frac{2}{\epsilon} \ln \frac{1}{\delta} \right\} \\ m &\geq \frac{2n^\alpha m^\beta \ln 2}{\epsilon} \\ m &\geq \left( \frac{2n^\alpha}{\epsilon} \ln 2 \right)^{\frac{1}{1-\beta}} \end{aligned}$$

□

מש"ל.

3.5.2 דוגמה - למידת OR של  $k$  משתנים

המטרה: להוריד את גודל המדגם מ- $O(n)$  ל- $O(k \log n)$ . לצורך הדוגמה, ניזכר בבעיית Set-Cover: הקלט הוא  $S_1, \dots, S_t \subseteq U$ , ואנחנו מחפשים  $S_{i_1}, \dots, S_{i_\ell} \subseteq U$  כך ש:

$$\bigcup_j S_{i_j} = U$$

נשתמש באלגוריתם החמדי לפתרון Set-Cover (אלגוריתם 2.3). ננתח אותו: נניח שיש כיסוי בגודל  $k$ . אזי:

$$\forall j. \exists t \in S_{\text{opt}}. |U_j \cap S_t| \geq \frac{|U_j|}{k}$$



**אלגוריתם 2.3 האלגוריתם החמדן ל-Set-Cover**1.  $U_0 \leftarrow U, j \leftarrow 0, S \leftarrow \emptyset$ 2. כל עוד  $U_j \neq \emptyset$ :(א) בחר  $S_i = \arg \max_{S_r} \{|S_i \cap U_j|\}$ (ב)  $S \leftarrow S \cup \{i\}$ (ג)  $U_{j+1} \leftarrow U_j \setminus S_i$ (ד)  $j \leftarrow j + 1$ 3. החזר  $S$ 

לכן:

$$\begin{aligned} |U_{j+1}| &\leq |U_j| - \frac{|U_j|}{k} = \left(1 - \frac{1}{k}\right) |U_j| \\ &= \left(1 - \frac{1}{k}\right)^{j+1} |U_0| \end{aligned}$$

עבור  $1 + k \ln |U|$  איטרציות, כיסינו את כל האיברים.  
 כדי ללמוד OR, נריך ELIM (אלגוריתם 1.3). נבצע רדוקציה ל-Set-Cover:

$$U = \{x \mid \langle x, + \rangle \in S\}$$

(כל הדגימות החיוביות)

$$S_{\ell_i} = \{x \in U \mid \ell_i \in x\}$$

נריך את האלגוריתם החמדן ל-Set-Cover (אלגוריתם 2.3), ונקבל  $1 + k \ln m$ .  
 גודל הקידוד:

$$\begin{aligned} (1 + k \ln m) \log(2n) &= O(l \ln m \ln n) \\ m &> \frac{k}{\varepsilon} \ln m \ln n + \frac{1}{\varepsilon} \ln \frac{1}{\delta} \end{aligned}$$

או בהצגה קצת שונה:

$$m > c \cdot \frac{k}{\varepsilon} \ln^2 n + \frac{1}{\varepsilon} \ln \frac{1}{\delta}$$

## פרק 4

# מודל Online<sup>1</sup>

נחשוב על הדוגמה הבאה: רובוט מסווג תפוזים לאיכות טובה או לא טובה. לאחר כל סיווג הוא מקבל משוב ממומחה. הרובוט יכול לעדכן את הפונקציה שלו, ואז לעבור לתפוז הבא. תיאור המודל הכללי שלנו:

1. האלגוריתם מקבל דגימה  $x$ .
2. האלגוריתם נותן תחזית לסיווג של  $x - b(x)$  (מכונה השערה נוכחית, או *current hypothesis*).
3. האלגוריתם מקבל את הסיווג הנכון -  $c^*(x)$  (זוהי פונקציית המטרה).
4. האלגוריתם ממשיך לדגימה הבאה.

נתבונן במודל *Adversarial* - כלומר, סדרת הקלט הגרועה ביותר שתיתן לנו את מספר השגיאות  $c^*(x) \neq b(x)$  המקסימלי.

### 4.1 למידה של מפריד לינארי

יש אוסף של נקודות חיוביות ואוסף של נקודות שליליות ב- $\mathbb{R}^n$  (או  $\{0, 1\}^n$ ). אנחנו רוצים למצוא וקטור  $w$  וחסם  $w_0$  כך ש- $w \cdot x = w_0$  (על מישור). האלגוריתם שלנו יהיה  $w \cdot x > 0$  גורר סיווג חיובי,  $w \cdot x < 0$  גורר סיווג שלילי. בלי הגבלת הכלליות,  $w_0 = 0$  (נוכל תמיד להוסיף קואורדינטה שתמיד תהיה 1).

#### 4.1.1 אלגוריתם Perceptron

הרעיון המרכזי העומד בבסיס האלגוריתם הוא שאם אין שגיאה, אין צורך לעדכן את ההשערה הנוכחית שלנו, ואם יש שגיאה, אז נעדכן את ההשערה הנוכחית בכיוון השגיאה. הנחת יסוד שלנו באלגוריתם תהיה ש- $\|x_t\| = 1$ , כלומר שכל הקלטים יהיו וקטור יחידה. אלגוריתם 1.4 מתאר את מהלך האלגוריתם Perceptron.

**משפט** תהי  $S$  סדרת דוגמאות מסווגות שעקבית עם מפריד לינארי  $w^*$  ( $\|w^*\| = 1$ ). אזי מספר השגיאות  $n$  של אלגוריתם Perceptron על הסדרה  $S$  חסום על ידי  $1/\gamma^2$ , כאשר:

$$\gamma = \min_{x \in S} \frac{|w^* \cdot x|}{\|x\|}$$

<sup>1</sup>שיעור שהתקיים בתאריך 11.11.2012. מבוסס על סיכומי של אולג.  
<sup>2</sup>לרוב זה לא ישפיע על האלגוריתם.

## אלגוריתם 1.4 Perceptron

1. נתחיל עם  $w_1 = \vec{0}$  ו- $t = 1$ .
2. בהינתן  $x_t$  נסווג + אם ורק אם  $w_t \cdot x_t > 0$ .
3. אם שגינו, נשים  $w_{t+1} \leftarrow w_t + \ell(x_t) \cdot x_t$ , כאשר  $\ell(x_t)$  הוא ה-label:

$$\ell(x_t) = \begin{cases} +1 & x_t \text{ is positive} \\ -1 & x_t \text{ is negative} \end{cases}$$

4. נמשיך ל- $t + 1$ .

$\gamma$  מכונה הפפריד או ה-margin. כאשר  $\|x\|$  מנורמל, זהו המרחק המינימלי מהעל-מישור.

**הוכחה** נתבונן בשני המדדים  $w_t \cdot w^*$  ו- $\|w_t\|$ . נניח כי בכל צעד האלגוריתם מבצע שגיאה (שאר המהלכים אינם רלוונטיים).

## טענת עזר 1

$$w_{t+1} \cdot w^* \geq w_t \cdot w^* + \gamma$$

**הוכחה** נניח כי  $x$  דגימה חיובית (ויש עליה שגיאה). אזי:

$$\begin{aligned} w_{t+1} \cdot w^* &= (w_t + x) \cdot w^* = w_t \cdot w^* + \underbrace{x \cdot w^*}_{>0} \\ &\geq w_t \cdot w^* + \gamma \end{aligned}$$

כאשר המעבר האחרון נובע מהגדרת  $\gamma$ .  
עבור  $x$  שלילי נקבל הוכחה דומה.

□

## טענת עזר 2

$$\|w_{t+1}\|^2 \leq \|w_t\|^2 + 1$$

**הוכחה** שוב נניח  $x$  דגימה חיובית. אזי:

$$\begin{aligned} \|w_{t+1}\|^2 &= \|w_t + x\|^2 = \|w_t\|^2 + \|x\|^2 + \underbrace{2x \cdot w_t}_{<0} \\ &\leq \|w_t\|^2 + \|x\|^2 = \|w_t\|^2 + 1 \end{aligned}$$

נשים לב כי  $2x \cdot w_t < 0$  כי הייתה שגיאה.

□

לאחר  $M$  שגיאות:

$$\begin{aligned} w_{M+1} \cdot w^* &\geq \gamma \cdot M \\ \|w_{M+1}\|^2 &\leq \sqrt{M} \end{aligned}$$

לכן:

$$\begin{aligned} \gamma \cdot M &\leq w_{M+1} \cdot w^* \leq w_{M+1} \cdot \frac{w_{M+1}}{\|w_{M+1}\|} \\ &= \|w_{M+1}\| \leq \sqrt{M} \\ M &\leq \frac{1}{\gamma^2} \end{aligned}$$

□

וסיימנו.

מה קורה אם אין מפריד מושלם?

נסמן:  $TD_\gamma$  - המרחק שצריך להזיז את הנקודות כדי לקבל מפריד של  $\gamma$ . טענת עזר 1 תהפוך ל:

$$w_{t+1} \cdot w^* \geq w_t \cdot w^* + \gamma - D_{\gamma,t}$$

כאשר  $D_{\gamma,t}$  הוא המרחק שצריך להזיז את  $x_t$  כדי לקבל את המפריד  $\gamma$ .  
לכן:

$$w_{M+1} \cdot w^* \geq \gamma \cdot M - TD_\gamma$$

טענת עזר 2 תישאר נכונה, ואז נקבל:

$$\sqrt{M} \geq \gamma \cdot M - TD_\gamma$$

חסם לפתרון:

$$M \leq \frac{1}{\gamma^2} + \frac{2}{\gamma} \cdot TD_\gamma$$

## 4.1.2 אלגוריתם Margin Perceptron

אלגוריתם 2.4 מתאר את מהלך האלגוריתם.

**משפט** לכל סדרת נקודות  $S$  עקביות עם מפריד  $w^*$  ( $\|w^*\| = 1$ ) ו- $\gamma = \min_{x \in S} \frac{w^* \cdot x}{\|x\|}$  אז מספר השגיאות חסום על ידי  $\frac{12}{\gamma^2}$ .

**הערה** הייתרון ב-Margin Perceptron הוא שמקביל מפריד יותר טוב עם margin גדול יותר.

## אלגוריתם 2.4 Margin Perceptron

1. נגדיר:  $w_1 = \ell(x_1) \cdot x_1$ .

2. נחזיר את התחזית שלנו:

(א) אם  $\frac{\gamma}{2} \leq \frac{w_t \cdot x}{\|w_t\|}$  אז נאמר חיובי.

(ב) אם  $-\frac{\gamma}{2} \geq \frac{w_t \cdot x}{\|w_t\|}$  אז נאמר שלילי.

(ג) אחרת נאמר שגיאה (margin mistake).

3. אם התחזית שלנו שגתה, נעדכן:

$$w_{t+1} \leftarrow w_t + \ell(x_t) \cdot x_t$$

$$t \leftarrow t + 1$$

## טענת עזר 1

$$w_{t+1} \cdot w^* \geq w_t \cdot w^* + \gamma$$

## הוכחה

$$\begin{aligned} w_{t+1} \cdot w^* &= (w_t + \ell(x_t) \cdot x_t) \cdot w^* \\ &= w_t \cdot w^* + \ell(x_t) \cdot x_t \cdot w^* \\ &\geq w_t \cdot w^* + \gamma \end{aligned}$$

## טענת עזר 2

$$\|w_{t+1}\| \leq \|w_t\| + \frac{1}{2 \cdot \|w_t\|} + \frac{\gamma}{2}$$

## 4.2 מודל Margin Bound

נניח כי  $c^* \in C$  פונקציית המטרה. בשלב  $t$ :

1. האלגוריתם מקבל  $x_t$ .

2. האלגוריתם בוחר סיווג  $b_t = h_t(x_t)$ .

3. האלגוריתם רואה את הסיווג הנכון  $c^*(x_t)$ .

**הגדרה** למחלקה  $C$  יש אלגוריתם  $A$  עם חסם שגיאה  $M$  אם לכל  $c^* \in C$  ולכל סדר דוגמאות  $S$  מספר השגיאות הוא לכל היותר  $M$ .

כמו כן, נניח כי  $C$  סופית.

**4.2.1 האלגוריתם (CON) Consistent**

בשלב  $t$  נגדיר את  $C_t$  להיות מחלקת כל ההשערות העקביות. נבחר  $h_t \in C_t$  ונחזיר  $b_t = h_t(x_t)$ .  
אזי:

$$C_{t+1} \subseteq C_t \bullet$$

$$\bullet \text{ אם יש שגיאה בזמן } t \text{ אז } C_{t+1} \subset C_t$$

$$\text{לכן, } M \leq |C| - 1$$

**4.2.2 אלגוריתם חציה (HAL)**

כמו קודם, נגדיר את  $C_t$  להיות מחלקת כל ההשערות העקביות בזמן  $t$ .  
כמו כן, נגדיר:

$$\text{one} = \{c \in C_t \mid c(x_t) = 1\}$$

$$\text{zero} = \{c \in C_t \mid c(x_t) = 0\}$$

התחזית תהיה 1 אם ורק אם  $|\text{one}| > |\text{zero}|$ . גם כאן:

$$C_{t+1} \subseteq C_t \bullet$$

$$\bullet \text{ אם יש שגיאה בשלב } t, \text{ אזי:}$$

$$|C_{t+1}| \leq \frac{|C_t|}{2}$$

$$\text{לכן נסיק כי } M \leq \log_2 |C|$$

**4.3 הקשר בין Mistake Bound ומודל PAC**

בהינתן אלגוריתם שמרני  $A$  עם חסם שגיאה  $M$ , אפשר להגדיר אלגוריתם PAC לאותה הבעיה,  $A_{\text{PAC}}$ . אלגוריתם 3.4 מציג כיצד ניתן לעשות זאת.

**משפט**  $A_{\text{PAC}}$  לומר PAC את  $C$ .

**הוכחה** בכל בלוק שלא עצרנו, ביצענו שגיאה אחת לפחות. כמו כן,  $A$  מבצע לכל היותר  $M$  שגיאות.

לכן, אם הגענו ל- $M$  שגיאות, ההשערה של  $A$  מושלמת. אחרת, בבלוק האחרון ראינו  $\frac{1}{\epsilon} \cdot \ln \frac{M}{\delta}$  דגימות, ואת כולן סיווגנו נכון. נניח של- $h_i$  יש שגיאה גדולה או שווה ל- $\epsilon$ . אזי:

$$\Pr [\forall j. h_i(x_j) = c^*(x_j)] \leq (1 - \epsilon)^{1/\epsilon \cdot \ln M/\delta} \leq \frac{\delta}{M}$$

כלומר, ההסתברות ש- $A_{\text{PAC}}$  יוציא השערה רעה (עם שגיאה גדולה או שווה ל- $\epsilon$ ) היא  $M \cdot \frac{\delta}{M} = \delta$ .  
□

**אלגוריתם 3.4** מציאת אלגוריתם PAC מ-Mistake Bound

1. ניקח מדגם בגודל  $M \cdot \frac{1}{\epsilon} \cdot \ln \frac{M}{\delta}$ .
  2. נחלק ל- $M$  קבוצות שוות.
  3. נריץ את  $A$  על הקבוצה ה- $i$ .
- (א) אם לא ביצע שגיאה, נחזיר את ההשערה הנוכחית.  
 (ב) אם ביצע שגיאה, נמשיך לקבוצה ה- $i + 1$ .
4. אם סיימנו את  $M$  הקבוצות, נחזיר את ההשערה הנוכחית.
- אם הגענו לשלב זה, סימן ש- $A$  ביצע  $M$  שגיאות. מכיוון שזהו חסם השגיאה, נקבל השערה מושלמת.

**אלגוריתם 4.4** Online למידה של OR

1. אתחול:  $L = \{x_1, \bar{x}_1, x_2, \bar{x}_2, \dots\}$ .
  2. בזמן  $t$ :
- (א) מקבלים דגימה  $z = z_1 z_2 \dots z_n$   
 (ב) נותנים תחזית לפי  $h_L(z)$   
 (ג) אם הייתה שגיאה:
- i. נגדיר  $S_z = \{\ell_i \mid \ell_i \text{ is positive in } z\}$ .
  - ii. נשים  $L \leftarrow L \setminus S_z$ .

**4.4 למידה של OR**

היה לנו אלגוריתם אלימינציה<sup>3</sup>. אלגוריתם 4.4 הוא גרסת Online שלו.

**משפט** מספר השגיאות יהיה לכל היותר  $n + 1$  שגיאות (כאשר  $|L|$  התחלתי הוא בגודל  $2n$ ).

**הוכחה** בשגיאה הרשונה, נפסול בדיוק  $n$  ליטרלים. כל שגיאה נוספת פוסלת לפחות ליטרל אחד נוסף.  $\square$

**4.5 אלגוריתם Winnow**

אלגוריתם 5.4 מתאר את אלגוריתם Winnow לחישוב מפריד לינארי. פה אנחנו פותרים את OR על ידי מפריד לינארי במרחב  $\{0, 1\}^n$ .

**משפט** אלגוריתם Winnow לומד OR של  $r$  משתנים חיוביים עם לכל היותר  $O(r \log n)$  שגיאות.

<sup>3</sup>ראה אלגוריתם 1.3.

**אלגוריתם 5.4 Winnow**

1. נאתחל  $w_0 = n, w^1 = (1, 1, \dots, 1)$  (בניגוד ל-0 שהיה קודם).

2. בהינתן נקודה  $x$ , נסווג אותה כחיובית אם ורק אם  $w^t \cdot x \geq n$ .

3. אם הייתה שגיאה:

(א) אם  $h(x) = 0$  ו- $c^*(x) = 1$ , אז  $w_i^{t+1} \leftarrow 2w_i^t$  ו- $x_i = 1$ .

(ב) אם  $h(x) = 1$  ו- $c^*(x) = 0$ , אז  $w_i^{t+1} \leftarrow \frac{w_i^t}{2}$  ו- $x_i = 1$ .

**הוכחה** נסמן:  $S = \{x_{i_1}, \dots, x_{i_r}\}$  כך ש- $c^*(x) = x_{i_1} \vee \dots \vee x_{i_r}$ . נסמן:  $w_r = \{w_{i_1}, \dots, w_{i_r}\}$ . נגדיר:  $w_i(t)$  - המשקל של  $x_i$  בזמן  $t$ . נסמן:

$$TW(t) = \sum_{i=1}^n w_i(t)$$

אם יש שגיאה על דגימה חיובית (כלומר  $c^*(x) = 1$ ), אז:

$$\exists i_j \in S. w_{i_j}(t+1) = 2w_{i_j}(t)$$

אם יש שגיאה על דגימה שלילית, אז:

$$\forall i_j \in S. w_{i_j}(t+1) = w_{i_j}(t)$$

לכן, לכל  $i_j \in S$ , היא פונקציה מונוטונית עולה. כמו כן, לכל  $i_j \in S$ , לא יכול לגדול יותר מ- $1 + \log n$  פעמים (כי אם הוא גדל כמות הזאת של פעמים, כבר לא נטעה אם  $i_j$  מופיע ב- $x$ ). לכן:

$$M_t \leq r \cdot (1 + \log n)$$

אם  $h(x) = 0$  ו- $c^*(x) = 1$ :

$$\begin{aligned} TW(t+1) &= TW(t) + \sum_{i=1}^n x_i w_i(t) \\ &\leq TW(t) + n \end{aligned}$$

אם  $h(x) = 1$  ו- $c^*(x) = 0$ :

$$\begin{aligned} TW(t+1) &= TW(t) - \frac{1}{2} \cdot \sum_{i=1}^n x_i w_i(t) \\ &\leq TW(t) - \frac{n}{2} \end{aligned}$$



נשים לב כי  $TW(t) > 0$  ואז:

$$0 < TW(t) \leq TW(0) + n \cdot M_+ - \frac{n}{2} \cdot M_-$$
$$M_- \leq 2M_+ + \frac{2}{n} \cdot \underbrace{TW(0)}_n = 2M_+ + 2$$

לכן נסיק כי:

$$M = M_- + M_+ \leq 3r \cdot (1 + \log n) + 2$$

□

וסיימנו.

## פרק 5

# <sup>1</sup>Regret Minimization

### 5.1 מבוא

לפי מודל Online, בזמן  $t$ :

1. מקבלים  $x_t$ .

2. נותנים תחזית  $b_t$ .

3. רואים את  $c^*(x_t)$ .

ראינו את אלגוריתם החציה (HAL): אם  $c^* \in H$ , אזי HAL יעשה לכל היותר  $O(\log |H|)$  טעויות.

אם  $c^* \notin H$ , נרצה לחזות כמו ההשערה הטובה ביותר ב- $H$ .

#### 5.1.1 המודל האלגוריתמי

המודל האלגוריתמי שלנו יהיה כזה: לכל  $h \in H$ , נשמור  $w_h$ , כך ש- $\sum_{h \in H} w_h = 1$  ו- $\forall h \in H. w_h \geq 0$ . התחזית תהיה ממוצע משוקלל. נסמן ב- $\ell_h^t$  את ההפסד של ההשערה  $h$  בזמן  $t$ . אזי  $\ell_h^t \in [0, 1]$ , ותוחלת ההפסד בזמן  $t$  היא:

$$\sum_{h \in H} \ell_h^t \cdot w_h$$

#### 5.1.2 External Regret

נניח כי:

$$L_h^T = \sum_{1 \leq t \leq T} \ell_h^t$$
$$L_{\text{best}}^T = \min_{h \in H} L_h^T$$

---

<sup>1</sup>שיעור שהתקיים בתאריך 18.11.2012.

**אלגוריתם 1.5 Deterministic Greedy**

1. עבור  $t = 1$  נבחר  $h_1$  (שרירותי).

2. עבור  $t > 1$  נבחר:

$$h_t = \arg \min_{h \in H} L_h^{t-1}$$

אזי ההפסד של האלגוריתם  $\mathcal{A}$  יהיה:

$$L_{\mathcal{A}}^T = \sum_{1 \leq t \leq T} \sum_{h \in H} \ell_h^t \cdot w_h^t$$

היעד שלנו יהיה למצוא אלגוריתם שיקיים את  $\text{Regret} \leq L_{\mathcal{A}}^T - L_{\text{best}}^T$ . המטרה שלנו היא שה- $\text{Regret}$  יהיה  $o(T)$ .

**5.2 אלגוריתמים****5.2.1 אלגוריתם Deterministic Greedy (G)**

אלגוריתם 1.5 מתאר את מהלך אלגוריתם Deterministic Greedy.

**משפט**

$$L_{\mathcal{G}}^T \leq |H| \cdot L_{\text{best}}^T + |H| - 1$$

**הוכחה** נגדיר את  $B_k$  להיות אוסף הזמנים בהם  $L_{\text{best}}^t = k$ . בצעד הראשון, ב- $B_k$  יש לכל היותר  $|H|$  השערות עבורן  $L_{\text{best}}^{t_0} = k$ . כל שגיאה תוריד את מספר ההשערות עם  $L_h = k$  בלפחות אחד. אחרי לכל היותר  $|H|$  שגיאות,  $L_{\text{best}}$  יגדל ל- $k+1$ . לכן:

$$L_{\mathcal{G}}^T \leq |H| \cdot L_{\text{best}}^T + |H| - 1$$

□ וסיימנו.

**משפט** לכל אלגוריתם דטרמיניסטי  $\mathcal{D}$  יש סדרת הפסדים עבורה:

$$L_{\mathcal{D}}^T \geq |H| \cdot L_{\text{best}}^T + (T \bmod |H|)$$

**הוכחה** נבחר סדרת הפסדים כך שבזמן  $t$ ,  $\mathcal{D}$  בוחר את  $h^t$  ומתקיים  $\ell_{h^t}^t = 1$ , ולכל  $h \neq h^t$  מתקיים  $\ell_h^t = 0$ . מבניית הסדרה,

$$L_{\mathcal{D}}^T = T$$

## אלגוריתם 2.5 Randomized Greedy

בזמן  $t$ , נסמן:

$$H^t = \left\{ h \in H \mid L_h^t = L_{\text{best}}^t \right\}$$

נגדיר התפלגות מעל  $H$  באמצעות:

$$p_h^t = \begin{cases} \frac{1}{|H^{t-1}|} & h \in H^{t-1} \\ 0 & \text{otherwise} \end{cases}$$

האלגוריתם יבחר  $h \in H$  לפי ההתפלגות שהגדרנו.

כמו כן:

$$\sum_{h \in H} L_h^T = T$$

לכן:

$$\exists h \in H. L_h^T \leq \left\lfloor \frac{T}{|H|} \right\rfloor$$

וסיימנו. □

## 5.2.2 אלגוריתם (GR) Randomized Greedy

אלגוריתם 2.5 מתאר את מהלך אלגוריתם Randomized Greedy.

משפט בתוחלת:

$$L_{\mathcal{RG}}^T \leq (\ln |H| + 1) \cdot L_{\text{best}}^T + \ln |H|$$

**הוכחה** נגדיר את  $B_k$  כמו קודם. היריב יעדיף לבצע שגיאה אחת בכל צעד. נניח שהוא עושה  $r$  שגיאות בבת אחת. אז ההפסד יהיה  $\frac{r}{m}$  (כאשר  $m = |H^{t-1}|$ ). אם הוא יעשה שגיאה אחת בכל פעם, נקבל:

$$\frac{1}{m} + \frac{1}{m-1} + \dots + \frac{1}{m-r+1}$$

נשים לב כי:

$$\frac{r}{m} < \frac{1}{m} + \frac{1}{m-1} + \dots + \frac{1}{m-r+1}$$

$$\begin{aligned} \mathbf{E} \left[ L_{\mathcal{RG}}^{B_k} \right] &= \frac{1}{m} + \frac{1}{m-1} + \dots + \frac{1}{m-r+1} \\ &= \sum_{i=1}^{|H|} \frac{1}{i} \leq \ln |H| + 1 \end{aligned}$$

**אלגוריתם 3.5** Randomized Weighted Majorityנגדיר את המשקל  $w_h$  של  $h \in H$  בזמן  $t$  להיות:

$$w_h^t = (1 - \eta)^{L_h^t}$$

נשים לב שאפשר לחשוב על הגדרה זו גם כעל הגדרה רקורסיבית:

$$w_h^{t+1} = \begin{cases} w_h^t & \ell_h^t = 0 \\ w_h^t \cdot (1 - \eta) & \ell_h^t = 1 \end{cases}$$

(כאשר  $w_h^0 = 1$ ).

נהפוך את המשקולות להסתברות: נגדיר:

$$W^t = \sum_{h \in H} w_h^t$$

נגדיר גם:

$$p_h^t = \frac{w_h^t}{W^t}$$

האלגוריתם יבחר השערה  $h \in H$  לפי ההתפלגות שהגדרנו.

לכן:

$$L_{\mathcal{RG}}^T \leq (\ln |H| + 1) \cdot L_{\text{best}}^T + \ln |H|$$

□

וסיימו.

**5.2.3 אלגוריתם Randomized Weighted Majority (RWM)**

אלגוריתם 3.5 מתאר את אלגוריתם Randomized Weighted Majority.

משפט עבור  $\eta \in (0, 1/2]$ :

$$L_{\mathcal{RWM}}^T \leq (1 + \eta) \cdot L_{\text{best}}^T + \frac{\log |H|}{\eta}$$

כמו כן, עבור  $\eta = \min \left\{ \frac{1}{2}, \sqrt{\frac{\log |H|}{T}} \right\}$ :

$$L_{\mathcal{RWM}}^T \leq L_{\text{best}}^T + 2 \cdot \sqrt{T \cdot \log |H|}$$

**הוכחה** נסמן ב- $F^t$  את ההסתברות לשגיאה של  $\mathcal{RWM}$ :

$$F^t = \frac{1}{W^t} \cdot \sum_{\substack{h \in H \\ \ell_h^t = 1}} w_h^t = \sum_{\substack{h \in H \\ \ell_h^t = 1}} p_h^t$$

אזי:

$$\begin{aligned}
 W^{t+1} &= \sum_{h \in H} w_h^{t+1} \\
 &= \sum_{\substack{h \in H \\ \ell_h^{t+1} = 0}} w_h^{t+1} + \sum_{\substack{h \in H \\ \ell_h^{t+1} = 1}} w_h^{t+1} \\
 &= \sum_{\substack{h \in H \\ \ell_h^{t+1} = 0}} w_h^t + \sum_{\substack{h \in H \\ \ell_h^{t+1} = 1}} w_h^t \cdot (1 - \eta) \\
 &= \sum_{h \in H} w_h^t - \eta \cdot \sum_{\substack{h \in H \\ \ell_h^{t+1} = 1}} w_h^t \\
 &= W^t - \eta \cdot F^t \cdot W^t \\
 &= W^t \cdot (1 - \eta \cdot F^t) \\
 &= W^1 \cdot \prod_{\tau=1}^t (1 - \eta \cdot F^\tau) \\
 &= |H| \cdot \prod_{\tau=1}^t (1 - \eta \cdot F^\tau)
 \end{aligned}$$

לכן:

$$\forall h \in H. W^{T+1} \geq w_h^{T+1} \geq (1 - \eta)^{L_{\text{best}}^T}$$

לכן:

$$(1 - \eta)^{L_{\text{best}}^T} \leq |H| \cdot \prod_{t=1}^T (1 - \eta \cdot F^t)$$

לכן:

$$L_{\text{best}}^T \cdot \ln(1 - \eta) \leq \ln |H| + \sum_{t=1}^T \ln(1 - \eta \cdot F^t)$$

נשתמש בזהות  $-z - z^2 \leq \ln(1 - z) \leq -z$

$$L_{\text{best}}^T (-1 - \eta^2) \leq \ln |H| + \sum_{t=1}^T (-\eta \cdot F^t)$$

לכן:

$$\eta \cdot \sum_{t=1}^T F^t \leq (\eta + \eta^2) \cdot L_{\text{best}}^T + \ln |H|$$

לכן:

$$L_{\mathcal{RWM}}^T \leq (1 + \eta) \cdot L_{\text{best}}^T + \frac{\log |H|}{\eta}$$

אם נבחר:

$$\eta = \min \left\{ \frac{1}{2}, \sqrt{\frac{\ln |H|}{T}} \right\}$$

אזי  $L_{\text{best}}^T \leq T$ , ואז:

$$L_{\mathcal{RWM}}^T \leq L_{\text{best}}^T + \eta \cdot T + \frac{\ln |H|}{\eta}$$

□

וסיימנו.

### 5.3 חסמים תחתונים לאלגוריתמי Online ממושקלים

אחרי שהראנו חסמים עליונים ל-External Regret, עלינו לבדוק כמה טוב יכול להיות אלגוריתם Online ממושקל. נבדוק את זה עבור שני מקרים.

#### 5.3.1 טווח קצר - $T = \frac{1}{2} \cdot \log |H|$

נבחר הפסדים כך ש:

$$\forall h \in H. \forall 1 \leq t \leq T. \Pr [\ell_h^t = 0] = \frac{1}{2} = \Pr [\ell_h^t = 1]$$

נראה שבהסתברות גבוהה קיים  $h \in H$  כך ש- $L_h^T = 0$ .  
ואכן:

$$\forall h \in H. \Pr [L_h^T = 0] = \left(\frac{1}{2}\right)^T = \frac{1}{\sqrt{|H|}}$$

לכן:

$$\Pr [L_h^T \neq 0] = 1 - \frac{1}{\sqrt{|H|}}$$

לכן:

$$\Pr [\forall h \in H. L_h^T \neq 0] = \left(1 - \frac{1}{\sqrt{|H|}}\right)^{|H|} \leq e^{-\sqrt{|H|}}$$

לכן:

$$\Pr [\exists h \in H. L_h^T = 0] \geq 1 - e^{-\sqrt{|H|}}$$

נעבור לכתיב תוחלות:

$$\begin{aligned} \mathbf{E} [L_{\text{best}}^T] &\leq T \cdot e^{-\sqrt{|H|}} = \frac{1}{2} \cdot \log |H| \cdot e^{-\sqrt{|H|}} \\ \mathbf{E} [L_{\text{ON}}] &= \frac{1}{2} \cdot T = \frac{1}{4} \cdot \log |H| \\ \mathbf{E} [\text{Regret}] &\geq \mathbf{E} [L_{\text{ON}}] - \mathbf{E} [L_{\text{best}}^T] \\ &\geq \frac{1}{4} \cdot \log |H| - \frac{1}{2} \cdot \log |H| \cdot e^{-\sqrt{|H|}} \\ &\approx \frac{1}{4} \cdot \log |H| \end{aligned}$$

### 5.3.2 כתלות בזמן - $|H| = 2$

נניח כי  $H = \{h_1, h_2\}$ . כמו כן, נניח כי סדרת ההפסדים היא  $(0, 1)$  בהסתברות  $1/2$ , ו- $(1, 0)$  בהסתברות  $1/2$ . אזי, אם האלגוריתם ON מחליט להחזיר  $h_1$  בהסתברות  $p$  ו- $h_2$  בהסתברות  $1-p$ :

$$\mathbf{E} [\ell_{\text{ON}}^t] = p \cdot \frac{1}{2} + (1-p) \cdot \frac{1}{2} = \frac{1}{2}$$

לכן:

$$\mathbf{E} [L_{\text{ON}}^T] = \frac{1}{2} \cdot T$$

ואז:

$$\begin{aligned} \mathbf{E} [L_{\text{best}}^T] &= \mathbf{E} [\min \{L_{h_1}^T, L_{h_2}^T\}] \\ &= \mathbf{E} \left[ \frac{T}{2} - \frac{|L_{h_1}^T - L_{h_2}^T|}{2} \right] \\ &= \frac{T}{2} - \frac{1}{2} \cdot \mathbf{E} [|L_{h_1}^T - L_{h_2}^T|] \\ &= \Omega(\sqrt{T}) \end{aligned}$$

לכן:

$$\mathbf{E} [\text{Regret}] \geq c \cdot \sqrt{T}$$

## Multi-Arm Bandit 5.4

המודל שלנו יהיה מעט שונה. בזמן  $t$ , כאשר מבצעים פעולה  $a \in A$ , מקבלים  $\ell_a^t$  (ולא את  $\ell_b^t$  אם  $a \neq b$ ). מודל סטוכסטי: נניח כי לכל פעולה יש מ"מ  $X_a$ , כך ש- $\mathbf{E}[\ell_a^t] = \mathbf{E}[X_a] = \mu_a$ . המטרה שלנו היא לבחור את  $a \in A$  כך ש- $\ell_a^t$  הכי נמוך.



## אלגוריתם 4.5 Test &amp; Play

1. שלב ה-Test: נדגום כל פעולה  $a \in A$  סדר גודל של  $O\left(\frac{1}{\varepsilon} \cdot \ln \frac{|A|}{\delta}\right)$  פעמים.
2. לכל פעולה  $a \in A$  נגדיר את ממוצע ההפסד  $\hat{\mu}_a$ .
3. שלב ה-Play: נבחר:

$$\hat{a}^* = \arg \min_{a \in A} \hat{\mu}_a$$

בשאר הפעולות, נבחר תמיד ב- $\hat{a}^*$ .

## 5.4.1 אלגוריתם Test &amp; Play

אלגוריתם Test & Play הוא האלגוריתם הפשוט ביותר ללמידת המודל. אלגוריתם 4.5 מתאר את מהלכו.

**משפט** בהסתברות  $1 - \delta$  נקבל:

$$\forall a \in A. |\mu_a - \hat{\mu}_a| \leq \frac{\varepsilon}{2}$$

□ **הוכחה** ישירות מחסמי Chernoff. כעת נרצה להעריך את  $|\mu^* - \mu_{\hat{a}^*}|$ . מהמשפט:

$$\mu_{\hat{a}^*} - \frac{\varepsilon}{2} \leq \hat{\mu}_{\hat{a}^*} \leq \hat{\mu}_{a^*} \leq \mu_{a^*} + \frac{\varepsilon}{2}$$

לכן:

$$|\mu^* - \mu_{\hat{a}^*}| \leq \frac{\varepsilon}{2}$$

נחשב חסם עליון על ה-Regret של Test & Play.

• בשלב ה-Test:

$$O\left(\frac{|A|}{\varepsilon^2} \cdot \ln \frac{|A|}{\delta}\right)$$

• בשלב ה-Play:

– בהסתברות  $1 - \delta$  בחרנו  $\varepsilon$ -טוב, ואז ה-Regret הוא  $\varepsilon \cdot T$ .

– בהסתברות  $\delta$ , ה-Regret הוא  $\delta \cdot T$ .

הסכום שלהם הוא תוחלת ה-Regret. נבחר  $\delta = \frac{1}{T}$ . נרצה להביא למינימום את:

$$\frac{|A|}{\varepsilon^2} \cdot \ln |A| \cdot T + \varepsilon \cdot T + 1$$

## אלגוריתם 5.5 Upper Confidence Bound

1. שלב האתחול: נדגום כל  $a \in A$  פעם אחת. מהדגימות האלה נקבל את  $\hat{\mu}_a$ , ו- $T_a = 1$ .

2. עבור  $t > |A|$ , נבחר את  $a^t$  שמוגדר על ידי:

$$a^t = \arg \min_{a \in A} \hat{\mu}_a - \sqrt{\frac{2 \cdot \log t}{T_a}}$$

כמו כן, נעדכן את  $\hat{\mu}_a$  ו- $T_a$  בהתאם.

זה נותן לנו  $\varepsilon = \frac{1}{T^{1/3}}$  בקירוב, או בדיוק:

$$\varepsilon = \left( \frac{2|A| \ln |A| \cdot T}{T} \right)^{1/3}$$

כלומר, בתנאים אופטימליים נבלה  $O(T^{2/3})$  ב-Test. ואז:

$$\mathbf{E}[\text{Regret}] \leq (2 \cdot |A| \ln |A| \cdot T)^{1/3} \cdot T^{2/3}$$

## 5.4.2 אלגוריתם (UCB) Upper Confidence Bound

אלגוריתם 5.5 מתאר את מהלכו של אלגוריתם Upper Confidence Bound.

**משפט** יהי  $\Delta_a = \mu_a - \mu^*$ . אזי:

$$\text{Regret (UCB)} \leq 8 \cdot \ln T \cdot \sum_{\substack{a \in A \\ a \neq a^*}} \frac{1}{\Delta_a} + 4 \cdot \sum_{\substack{a \in A \\ a \neq a^*}} \frac{1}{\Delta_a}$$

**הוכחה** ניזכר כי:

$$\mathbf{E}[\text{Regret}] = \sum_{\substack{a \in A \\ a \neq a^*}} \Delta_a \cdot \mathbf{E}[T_a]$$

לכן, נרצה לחסום את  $\mathbf{E}[T_a]$ :

$$\begin{aligned} T_a &= 1 + \sum_{t=|A|+1}^T I\{a^t = a\} \\ &\leq \ell + \sum_{t=|A|+1}^T I\{a^t = a, T_a \geq \ell\} \\ &\leq \ell + \sum_t \sum_{\ell \leq r \leq T} \sum_{1 \leq s \leq T} I\left\{ \hat{\mu}_a - \sqrt{\frac{2 \log T}{r}} \leq \hat{\mu}_{a^*}^s - \sqrt{\frac{2 \log t}{s}} \right\} \end{aligned}$$

האינדיקטור מתקיים בלפחות אחד מהתנאים הבאים:

$$\hat{\mu}_{a^*}^s \geq \mu^* + \sqrt{\frac{2 \cdot \log t}{s}} \quad (5.4.1)$$

$$\mu_a^r \leq \mu_a - \sqrt{\frac{2 \cdot \log t}{r}} \quad (5.4.2)$$

$$\mu^* > \mu_a + 2 \cdot \sqrt{\frac{2 \cdot \log t}{r}} \quad (5.4.3)$$

את 5.4.1 ו-5.4.2 נחסום עם Chernoff בהסתברות  $\frac{1}{t^4}$ . לגבי 5.4.3: אם  $\ell = \left\lceil \frac{8 \cdot \log t}{\Delta^2} \right\rceil$ , אז 5.4.3 לא מתקיים, ואז נישאר עם  $T_a \leq \ell + \sum \frac{1}{t^2}$  ו- $\sum \frac{1}{t^2}$  קבוע.  $\square$

## פרק 6

# Boosting<sup>1</sup>

### 6.1 למידה חלשה וחזקה

במודל PAC, מקבלים דוגמאות  $\langle x, c^*(x) \rangle$ , וקיימת התפלגות  $D$  עבור  $x$ . המטרה היא למצוא השערה  $h \in H$  כך ש- $\text{error}(h, c^*) \leq \varepsilon$  בהסתברות  $1 - \delta$ .  $\delta$  נקרא פרמטר הבטחון, ו- $\varepsilon$  נקרא פרמטר הדיוק. שני הפרמטרים ביחד קובעים את חוזק הלמידה שלנו.

#### 6.1.1 שיפור בפרמטר הבטחון

ההנחה: נתון אלגוריתם  $A$  שבהסתברות לפחות  $1/2$  מחזיר השערה  $h \in H$  כך ש:

$$\text{error}(h, c^*) \leq \varepsilon$$

אלגוריתם 1.6 מתאר את אלגוריתם Boost-Confidence, שמשפר את פרמטר הבטחון. ננתח את האלגוריתם:

1. לכל  $i$ :

$$\Pr \left[ \text{error}(h_i) \leq \frac{\varepsilon}{2} \right] \geq \frac{1}{2}$$

לכן:

$$\Pr \left[ \exists i. \text{error}(h_i) \leq \frac{\varepsilon}{2} \right] \geq 1 - \left( \frac{1}{2} \right)^k$$

עבור  $k = \log \frac{2}{\delta}$ :

$$\begin{aligned} \Pr \left[ \exists i. \text{error}(h_i) \leq \frac{\varepsilon}{2} \right] &\geq 1 - \left( \frac{1}{2} \right)^k \\ &= 1 - \left( \frac{1}{2} \right)^{\log \frac{2}{\delta}} \\ &= 1 - \frac{1}{2^{\log \frac{2}{\delta}}} \\ &= 1 - \frac{1}{2/\delta} = 1 - \frac{\delta}{2} \end{aligned}$$

---

<sup>1</sup>שיעור שהתקיים בתאריך 25.11.2012.

**אלגוריתם 1.6 Boost-Confidence**

נגדיר את אלגוריתם Boost-Confidence (BC) שפועל על האלגוריתם  $\mathcal{A}$  (אותו הוא מקבל כקלט):

1. הרץ את האלגוריתם  $\mathcal{A}$  במשך  $k = \log \frac{2}{\delta}$  פעמים עם רמת דיוק  $\varepsilon/2$ . הפלט שנקבל הוא  $h_1, \dots, h_k$ .

2. נקח דגימה נוספת  $S$  בגודל:

$$m = \frac{2}{\varepsilon^2} \cdot \ln \frac{4k}{\delta} = O\left(\frac{1}{\varepsilon^2} \cdot \ln \frac{k}{\delta}\right)$$

עבור כל  $h_i$ , נחשב את  $\widehat{\text{error}}(h_i)$ .

3. נחזיר את:

$$\hat{h}^* = \arg \min_{h_i} \widehat{\text{error}}(h_i)$$

2. עבור כל  $h_i$ : ניזכר כי:

$$\widehat{\text{error}}(h_i, c^*) = \frac{1}{m} \cdot \sum_{j=1}^m I(h_i(x_j) \neq c^*(x_j))$$

לכן, מ-Chernoff:

$$\Pr\left[|\widehat{\text{error}}(h_i, c^*) - \text{error}(h_i, c^*)| > \frac{\varepsilon}{2}\right] \leq 2 \cdot e^{-2 \cdot \left(\frac{\varepsilon}{2}\right)^2 \cdot m}$$

נרצה לחסום את ההסתברות למאורע ה"רע" הזה על ידי  $\frac{\delta}{2}$ , ולכן נדרוש:

$$2 \cdot k \cdot e^{-2 \cdot \left(\frac{\varepsilon}{2}\right)^2 \cdot m} \leq \frac{\delta}{2}$$

נפתור את המשוואה:

$$\begin{aligned} 2 \cdot k \cdot e^{-2 \cdot \left(\frac{\varepsilon}{2}\right)^2 \cdot m} &\leq \frac{\delta}{2} \\ 4 \cdot k \cdot e^{-2 \cdot \left(\frac{\varepsilon}{2}\right)^2 \cdot m} &\leq \delta \\ \ln 4k - 2 \cdot \left(\frac{\varepsilon}{2}\right)^2 \cdot m &\leq \ln \delta \\ 2 \cdot \left(\frac{\varepsilon}{2}\right)^2 \cdot m &\geq \ln 4k - \ln \delta \\ \frac{\varepsilon^2}{2} \cdot m &\geq \ln \frac{4 \cdot k}{\delta} \\ m &\geq \frac{2}{\varepsilon^2} \cdot \ln \frac{4 \cdot k}{\delta} \end{aligned}$$

**אלגוריתם 2.6** אלגוריתם לשקילת למידה חלשה וחזקה

הקלט: דגימות  $x_1, \dots, x_m$  וסיווגן על ידי  $c^*$ .  
 נניח כי  $H$  היא מחלקת השערות (חלשות), ו- $\mathcal{RM}$  הוא אלגוריתם Regret-Minimization.

1. בכל שלב  $t$ , האלגוריתם  $\mathcal{RM}$  יבחר פילוג  $D_t$  מעל  $x_1, \dots, x_m$ .
2. בהינתן  $D_t$ , קיימת השערה  $h_t \in H$  עבורה:

$$\text{error}_{D_t}(h_t) \leq \frac{1}{2} - \gamma$$

3. ההפסד יהיה 1 לכל סיווג נכון של  $h_t$ , ו-0 לכל סיווג מוטעה.
4. אחרי  $T$  שלבים, נחזיר את  $\text{MAJ}(h_1, \dots, h_T)$ .

לכן:

$$\begin{aligned} \text{error}(\hat{h}^*) &\leq \widehat{\text{error}}(\hat{h}^*) + \frac{\varepsilon}{2} \\ &\leq \widehat{\text{error}}(h_i) + \frac{\varepsilon}{2} \\ &\leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon \end{aligned}$$

**6.1.2 שיפור בפרמטר הדיוק**

נראה שקיים  $\gamma > 0$  כך שלכל פונקציית מטרה  $c^* \in C$ , לכל התפלגות  $D$  ולכל פרמטר דיוק  $\delta$ , אלגוריתם  $\mathcal{A}$  מוצא השערה  $h \in H$  כך ש:

$$\text{error}(h, c^*) \leq \frac{1}{2} - \gamma$$

**דוגמה** אם  $x_1 = x_2 = 1$ , אז  $c^*(x)$  היא פונקציה קשה. אחרת,  $c^*(x) = 0$ . נניח כי  $D$  היא ההתפלגות האחידה.

ניתן להגיע ל-87.5% דיוק בקלות. כדי להעלות את החוזי הדיוק, היינו משנים את  $D$  כך שנתמקד במקרה  $x_1 = x_2 = 1$ . זאת הסיבה לכך שאנחנו נראה שקיים  $\gamma$  לכל  $D$ , ולא שלכל  $D$  קיים  $\gamma$ .

**הוכחת שקילות למידה חלשה וחזקה** אלגוריתם 2.6 הוא האלגוריתם שיעזור לנו להוכיח את השקילות. ננתח אותו:

1. בכל שלב מצאנו  $h_t$ . ההפסד של  $\mathcal{RM}$  הוא  $(\frac{1}{2} + \gamma)$ . לכן, סך כל ההפסד יהיה  $(\frac{1}{2} + \gamma) \cdot T$ .

2. אם קיים  $x_i$  כך ש- $\text{MAJ}(h_1, \dots, h_T)$  טועה עליו, אזי ההפסד של  $x_i$  הוא לכל היותר  $\frac{T}{2}$ . ה- $\text{Regret}$  יהיה:

$$\left(\frac{1}{2} + \gamma\right) \cdot T \leq \frac{T}{2} + 2 \cdot \sqrt{T \cdot \log m}$$

נחלץ את  $T$ :

$$\begin{aligned}\gamma \cdot T &\leq 2 \cdot \sqrt{T \cdot \log m} \\ \gamma^2 \cdot T^2 &\leq 4 \cdot T \cdot \log m \\ \gamma^2 \cdot T &\leq 4 \cdot \log m \\ T &\leq \frac{4 \cdot \log m}{\gamma^2}\end{aligned}$$

## 6.2 בנייה רקורסיבית

נמצא שלוש השערות  $h_1, h_2, h_3$  כד ש- $\text{MAJ}(h_1, h_2, h_3)$  משפר את הדיוק. אלגוריתם 3.6 מתאר איך עושים את זה.

**הערכת השגיאה** נעריך את השגיאה של האלגוריתם. מה תהיה השגיאה אם ההסתברות של  $h_1, h_2, h_3$  ב"ת? כלומר, מה תהיה השגיאה אם ההסתברות לשגיאה של  $h_1, h_2, h_3$  ב"ת?

$$\text{error} = p^3 + 3p^2(1-p) = 3p^2 - 2p^3$$

ברור שההסתברויות האלה אינן ב"ת, מכיוון ש- $h_2$  נבנתה על סמך  $h_1$ , ו- $h_3$  נבנתה בהתבסס על  $h_1$  ו- $h_2$ . למרות שההסתברויות אינן ב"ת, גם במודל שלנו נקבל שגיאה כזו. נחלק את הדוגמאות לארבע קבוצות:

$$\begin{aligned}S_{cc} &= \{x \mid h_1(x) = c^*(x) = h_2(x)\} \\ S_{ce} &= \{x \mid h_1(x) = c^*(x) \neq h_2(x)\} \\ S_{ec} &= \{x \mid h_1(x) \neq c^*(x) = h_2(x)\} \\ S_{ee} &= \{x \mid h_1(x) \neq c^*(x) \neq h_2(x)\}\end{aligned}$$

ההסתברויות של הקבוצות האלו ביחס ל- $D_1$  הן  $p_{cc}, p_{ce}, p_{ec}$  ו- $p_{ee}$  בהתאמה. השגיאה הכוללת שלנו תהיה:

$$\text{error} = p_{ee} + (p_{ce} + p_{ec}) \cdot p$$

(נזכיר כי  $p$  היא הסתברות השגיאה של  $h_3$ ).  
נסמן:

$$\alpha = D_2(S_{ce})$$

לפי הגדרת  $D_2$ , ביחס ל- $D_1$ :

$$p_{ce} = 2 \cdot (1-p) \cdot \alpha$$

כמו כן:

$$D_2(S_{ec}) = p - \alpha$$

כי:

$$p = D_2(S_{ce}) + D_2(S_{ec})$$

**אלגוריתם 3.6** אלגוריתם לשיפור הדיוק

נניח כי  $\mathcal{A}$  הוא אלגוריתם ללמידה חלשה, ו- $p$  תהיה הסתברות השגיאה של  $\mathcal{A}$ . כמו כן, נניח כי נתון לנו Oracle מעל מרחב הדגימות, שנותן לנו דגימה לפי ההתפלגות  $D$ .

1. נסמן:  $D_1 = D$ . נמצא השערה  $h_1 \in H$  ביחס ל- $D_1$  (לפי  $\mathcal{A}$ ). השגיאה חסומה על ידי  $\frac{1}{2} - \gamma = p$ .

2. נגדיר:

$$\begin{aligned} S_c &= \{x \mid h_1(x) = c^*(x)\} \\ S_e &= \{x \mid h_1(x) \neq c^*(x)\} \end{aligned}$$

3. נגדיר התפלגות חדשה  $D_2$  כך ש- $D_2(S_e) = \frac{1}{2}$  וחוץ מזה  $D_2$  תהיה בדיוק כמו  $D_1$ . כלומר:

$$D_2 = \begin{cases} \frac{0.5}{1-p} \cdot D_1(x) & x \in S_c \\ \frac{0.5}{p} \cdot D_1(x) & x \in S_e \end{cases}$$

איך נייצור את  $D_2$  בפועל? נתון לנו Oracle שמייצר  $D_1$ . אנחנו רוצים להביא ל- $1/2$  את כמות הפעמים ש- $h_1$  צודקת ואת כמות הפעמים ש- $h_1$  טועה. לכן, באופן הזה נבנה Oracle שפועל לפי  $D_2$ , כי נחליט (בהתפלגות אחידה) אם הפעם אנחנו רוצים ש- $h_1$  תצדק או ש- $h_1$  תטעה.

4. נרוץ עם  $\mathcal{A}$  עם ההתפלגות  $D_2$  ונמצא  $h_2$  (עם שגיאה  $\frac{1}{2} - \gamma = p$ ).

5. נגדיר התפלגות  $D_3$  שתשים משקל על  $\{x \mid h_1(x) \neq h_2(x)\}$ . כלומר:

$$D_3(x) = \begin{cases} \frac{D_1(x)}{z} & h_1(x) \neq h_2(x) \\ 0 & h_1(x) = h_2(x) \end{cases}$$

כאשר  $z = \Pr[h_1(x) \neq h_2(x)]$  הוא פוקטור נרמול.

קל לבנות Oracle שפועל לפי  $D_3$ , באופן דומה לבניית ה-Oracle שפועל לפי  $D_2$ .

6. מוצאים עם  $\mathcal{A}$  השערה  $h_3$  ביחס ל- $D_3$  (עם שגיאה  $\frac{1}{2} - \gamma = p$ ).

7. הפלט שלנו יהיה:

$$H(x) = \begin{cases} h_1(x) & h_1(x) = h_2(x) \\ h_3(x) & h_1(x) \neq h_2(x) \end{cases}$$

או במילים אחרות,  $H = \text{MAJ}(h_1, h_2, h_3)$ .



לכן:

$$p_{ee} = 2 \cdot p \cdot (p - \alpha)$$

באופן דומה, נקבל כי:

$$D_2(S_{ec}) = \frac{1}{2} - (p - \alpha)$$

ולכן:

$$p_{ec} = 2 \cdot p \cdot \left( \frac{1}{2} - p + \alpha \right)$$

נציב ב-error:

$$\begin{aligned} \text{error} &= p_{ee} + (p_{ce} + p_{ec}) \cdot p \\ &= 2 \cdot p \cdot (p - \alpha) + \left( 2 \cdot p \cdot \left( \frac{1}{2} - p + \alpha \right) + 2 \cdot (1 - p) \cdot \alpha \right) \cdot p \\ &= 2p^2 - \cancel{2 \cdot p \cdot \alpha} + p^2 - 2p^3 + \cancel{2 \cdot p^2 \cdot \alpha} + \cancel{2 \cdot p \cdot \alpha} - \cancel{2 \cdot p^2 \cdot \alpha} \\ &= 3p^2 - 2p^3 \end{aligned}$$

למעשה, היינו רוצים להראות ש- $\gamma$  הולך וגדל עם הזמן. כלומר, אם

$$p_0 = \frac{1}{2} - \gamma_0$$

אזי

$$\frac{1}{2} - \gamma_{t+1} = 3 \cdot \left( \frac{1}{2} - \gamma_t \right)^2 - 2 \cdot \left( \frac{1}{2} - \gamma_t \right)^3$$

ולכן:

$$\frac{1}{2} - \gamma_{t+1} \leq \frac{1}{2} - \gamma_t \cdot \left( \frac{3}{2} - \gamma_t \right)$$

דרך אחרת להסתכל על זה היא באמצעות:

$$p_{t+1} \leq 3 \cdot p_t^2 - 2 \cdot p_t^3$$

לכן, אם  $\gamma$  היא שגיאת ההשערה החלשה, ו- $\epsilon$  היא שגיאת ההשערה החזקה, עומק הרקורסיה שלנו יהיה  $O\left(\log \log \frac{1}{\epsilon} + \log \frac{1}{\gamma}\right)$ . גודל הרקורסיה שלנו יהיה  $3^{\text{depth}}$ .

## 6.3 אלגוריתם AdaBoost

אלגוריתם 4.6 הוא אלגוריתם AdaBoost (Adaptive Boost).

**אלגוריתם 4.6 (Adaptive Boost) AdaBoost**

קלט: נקודות  $x_1, \dots, x_m$  מסווגות  $\langle x_i, y_i \rangle$ .  
 נגדיר:  $D_t$  - ההתפלגות בשלב  $t$ .  $D_t(i)$  - המשקל של הנקודה  $x_i$  בשלב  $t$ .

1. אתחול:

$$\forall i. D_1(i) = \frac{1}{m}$$

2. בשלב  $t$ :

(א) נבחר  $h_t \in H$  שמביאה למינימום את השגיאה  $\varepsilon_t$  ביחס ל- $D_t$ :

$$\varepsilon_t = \Pr_{D_t}[h_t(x) \neq c^*(x)]$$

(ב) נגדיר את  $D_{t+1}$ . נניח כי  $h_t(x) \in \{\pm 1\}$  אזי:

$$D_{t+1}(i) = \frac{D_t(i)}{z_t} \cdot \begin{cases} e^{-\alpha_t} & h_t(x_i) = y_i \\ e^{+\alpha_t} & h_t(x_i) \neq y_i \end{cases}$$

כאשר:

$$\alpha_t = \frac{1}{2} \cdot \ln \frac{1 - \varepsilon_t}{\varepsilon_t}$$

ו- $z_t$  הוא גורם מנרמל.

3. אחרי  $T$  שלבים:

$$H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t \cdot h_t(x) \right)$$

**משפט**

$$\begin{aligned} \widehat{\text{error}}(H) &\leq \prod_{t=1}^T \sqrt{4 \cdot \varepsilon_t \cdot (1 - \varepsilon_t)} \\ &= \prod_{t=1}^T \sqrt{4 \cdot \left(\frac{1}{2} - \gamma_t\right) \cdot \gamma_t} \end{aligned} \quad (6.3.1)$$

$$\begin{aligned} &= \prod_{t=1}^T \sqrt{1 - 4 \cdot \gamma_t^2} \\ &= \prod_{t=1}^T (1 - 4 \cdot \gamma_t^2)^{1/2} \\ &\leq \prod_{t=1}^T (e^{-4 \cdot \gamma_t^2})^{1/2} \\ &= \prod_{t=1}^T e^{-2 \cdot \gamma_t^2} \\ &= e^{-2 \cdot \sum_{t=1}^T \gamma_t^2} \end{aligned} \quad (6.3.2)$$

נסביר את המעברים בשלבים:

- מעבר 6.3.1 בוצע כי  $\varepsilon_t = 1 - \gamma_t$ .
- מעבר 6.3.2 בוצע לפי הזהות  $1 + x \leq e^x$ .

**הוכחה** שלג א': נוכיח כי:

$$D_{T+1}(i) = \frac{D_1(i) \cdot e^{-y_i \cdot f(x_i)}}{\prod_{t=1}^T z_t}$$

כאשר:

$$f(x) = \sum_{t=1}^T \alpha_t \cdot h_t(x)$$

מהרקורסיה נקבל:

$$D_{t+1}(i) = \frac{D_t(i)}{z_t} \cdot e^{-y_i \cdot \alpha_t \cdot h_t(x_i)}$$

לכן:

$$\begin{aligned} D_{T+1}(i) &= \frac{D_1(i)}{\prod_{t=1}^T z_t} \cdot e^{-y_i \cdot \sum_{t=1}^T \alpha_t \cdot h_t(x_i)} \\ &= \frac{D_1(i)}{\prod_{t=1}^T z_t} \cdot e^{-y_i \cdot f(x_i)} \\ &= \frac{e^{-y_i \cdot f(x_i)}}{m \cdot \prod_{t=1}^T z_t} \end{aligned} \quad (6.3.3)$$

כאשר מעבר 6.3.3 בוצע על סמך העובדה:

$$D_1(i) = \frac{1}{m}$$

שלג ב! נוכיח כי:

$$\widehat{\text{error}}(H) \leq \prod_{t=1}^T z_t$$

נחזור להגדרה של  $\widehat{\text{error}}(H)$ :

$$\begin{aligned} \widehat{\text{error}}(H) &= \frac{1}{m} \cdot \sum_{i=1}^m I(y_i \cdot f(x_i) \leq 0) \\ &\leq \frac{1}{m} \cdot \sum_{i=1}^m e^{-y_i \cdot f(x_i)} \\ &= \frac{1}{\mathcal{M}} \cdot \sum_{i=1}^m \mathcal{M} \cdot D_{T+1}(i) \cdot \prod_{t=1}^T z_t \end{aligned} \quad (6.3.4)$$

$$\begin{aligned} &= \prod_{t=1}^T z_t \cdot \sum_{i=1}^m D_{T+1}(i) \\ &= \prod_{t=1}^T z_t \end{aligned} \quad (6.3.5)$$

כאשר המעבר 6.3.4 בוצע על סמך הטענה מחלק א', והמעבר 6.3.5 בוצע על סמך העובדה:

$$\sum_{i=1}^m D_t(i) = 1$$

כי  $D_t$  מתאר התפלגות (עם נרמול). שלג ג! נוכיח כי:

$$z_t = 2 \cdot \sqrt{\varepsilon_t \cdot (1 - \varepsilon_t)}$$

$z_t$  הוא הגורם המנרמל של  $D_t$ . כלומר, לפי ההגדרה:

$$\begin{aligned} z_t &= \sum_{i=1}^m D_t(i) \cdot e^{-y_i \cdot \alpha_t \cdot h_t(x_i)} \\ &= \sum_{\substack{i=1 \\ h_t(x_i) = y_i}}^m D_t(i) \cdot e^{-\alpha_t} + \sum_{\substack{i=1 \\ h_t(x_i) \neq y_i}}^m D_t(i) \cdot e^{\alpha_t} \\ &= (1 - \varepsilon_t) \cdot e^{-\alpha_t} + \varepsilon_t \cdot e^{\alpha_t} \end{aligned}$$

נחשוב על  $z_t$  כפונקציה של  $\alpha_t$ , ונגזור:

$$\begin{aligned} \frac{\partial}{\partial \alpha_t} z_t &= -(1 - \varepsilon_t) \cdot e^{-\alpha_t} + \varepsilon_t \cdot e^{\alpha_t} = 0 \\ -e^{-\alpha_t} + \varepsilon_t \cdot e^{-\alpha_t} + \varepsilon_t \cdot e^{\alpha_t} &= 0 \\ -1 + \varepsilon_t + \varepsilon_t \cdot e^{2 \cdot \alpha_t} &= 0 \\ \varepsilon_t \cdot e^{2 \cdot \alpha_t} &= 1 - \varepsilon_t \\ e^{2 \cdot \alpha_t} &= \frac{1 - \varepsilon_t}{\varepsilon_t} \\ 2 \cdot \alpha_t &= \ln \frac{1 - \varepsilon_t}{\varepsilon_t} \\ \alpha_t &= \frac{1}{2} \cdot \ln \frac{1 - \varepsilon_t}{\varepsilon_t} \end{aligned}$$

אם נציב את  $\alpha_t$  שקיבלנו, נקבל:

$$\begin{aligned} z_t &= (1 - \varepsilon_t) \cdot e^{-\alpha_t} + \varepsilon_t \cdot e^{\alpha_t} \\ &= (1 - \varepsilon_t) \cdot e^{-\frac{1}{2} \cdot \ln \frac{1 - \varepsilon_t}{\varepsilon_t}} + \varepsilon_t \cdot e^{\frac{1}{2} \cdot \ln \frac{1 - \varepsilon_t}{\varepsilon_t}} \\ &= (1 - \varepsilon_t) \cdot \left( \frac{1 - \varepsilon_t}{\varepsilon_t} \right)^{-1/2} + \varepsilon_t \cdot \left( \frac{1 - \varepsilon_t}{\varepsilon_t} \right)^{1/2} \\ &= (1 - \varepsilon_t) \cdot \left( \frac{\varepsilon_t}{1 - \varepsilon_t} \right)^{1/2} + \varepsilon_t \cdot \left( \frac{1 - \varepsilon_t}{\varepsilon_t} \right)^{1/2} \\ &= (1 - \varepsilon_t) \cdot \frac{\varepsilon_t^{1/2}}{(1 - \varepsilon_t)^{1/2}} + \varepsilon_t \cdot \frac{(1 - \varepsilon_t)^{1/2}}{\varepsilon_t^{1/2}} \\ &= \varepsilon_t^{1/2} \cdot (1 - \varepsilon_t)^{1/2} + \varepsilon_t^{1/2} \cdot (1 - \varepsilon_t)^{1/2} \\ &= \sqrt{\varepsilon_t \cdot (1 - \varepsilon_t)} + \sqrt{\varepsilon_t \cdot (1 - \varepsilon_t)} \\ &= 2 \cdot \sqrt{\varepsilon_t \cdot (1 - \varepsilon_t)} \end{aligned}$$

□

מהחיבור של שלושת השלבים הוכחנו את המשפט.

## פרק 7

# <sup>1</sup>Nearest Neighbor

### 7.1 הקדמה

נתון מרחב מדגם  $X$ . המטרה שלנו היא לסווג נקודה  $x \in X$ . כמו כן, נתון מגדם  $S = \{x_i, b_i\}$ . אנחנו נמצא נקודה  $y \in S$  כך ש- $\|x - y\| = \min$  ואם  $\langle y, b \rangle$  אז ההשערה עבור  $x$  היא  $b$ .

### 7.2 שיטות כלליות

יש לנו כמה שיטות כלליות:

1. הנקודה הקרובה ביותר.

2. ממוצע של  $k$  נקודות קרובות ביותר.

3. משקל על הנקודות לפי מרחק.

### 7.3 מודלים ל-Nearest Neighbor

כמו כן, ישנם מספר מודלים לייצור הנקודות:

- נניח שישנה התפלגות  $D$  מעל  $X$ . תחילה דוגמים מ- $X$  לפי  $D$ . בהינתן  $x \in X$ , ישנה פונקציה  $p(x)$  כך ש:

$$\Pr[b = 1 | x] = p(x)$$

- נניח שישנן שתי התפלגויות  $D_0$  ו- $D_1$ , ופרמטר  $\alpha \in [0, 1]$ . תחילה נדגום את  $b$  לפי  $\alpha$  (כלומר  $\Pr[b = 1] = \alpha$ ), ואז נדגום את  $x \in X$  לפי  $D_b$ . מחזירים את  $\langle x, b \rangle$ .

$$p(x) = \frac{\alpha \cdot D_1(x)}{\alpha \cdot D_1(x) + (1 - \alpha) \cdot D_0(x)}$$

---

<sup>1</sup>שיעור שהתקיים בתאריך 1.12.2012.

- נניח שקיימת התפלגות  $D$  מעל  $X \times \{0, 1\}$  כשנדגום את  $D$  נקבל  $\langle x, b \rangle$  ישירות.

$$p(x) = \frac{D(\langle x, 1 \rangle)}{D(\langle x, 1 \rangle) + D(\langle x, 0 \rangle)}$$

אנחנו רוצים להביא למינימום את מספר השגיאות.

### 7.3.1 0-1 Loss

$$h(x) = \begin{cases} 1 & p(x) \geq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

### 7.3.2 Bayes Risk

ההפסד של המסווג האופטימלי (בהינתן ההתפלגות):

$$r(x) = \min \{p(x), 1 - p(x)\}$$

ואז:

$$R^* = \mathbf{E}_x [\text{Loss}(h, x)] = \int_x r(x) \cdot D(x) dx$$

נסמן:  $P = \mathbf{E}_x [\text{Loss}(\text{NN}, x)]$  (Nearest Neighbor - NN). נוכיח כי:

$$R^* \leq P \leq 2 \cdot R^* \cdot (1 - R^*)$$

אם  $R^* \sim \varepsilon$ , אזי צ"ל  $P \leq 2 \cdot \varepsilon \cdot (1 - \varepsilon)$ . אם  $R^* \sim \frac{1}{2} - \varepsilon$ , אזי  $P \leq 2 \cdot (\frac{1}{2} - \varepsilon) \cdot (\frac{1}{2} + \varepsilon) = \frac{1}{2} - 2 \cdot \varepsilon^2$ .

#### 7.3.2.1 מקרה פשוט

רוצים לסווג את  $x \in X$ . במדגם יש נקודה אחת  $\langle x_i, b_i \rangle$  כך ש- $x_i \equiv x$ . נחשב:

$$\begin{aligned} \Pr[b_i \neq b] &= 2 \cdot p(x) \cdot (1 - p(x)) \\ &= 2 \cdot r(x) \cdot (1 - r(x)) \end{aligned}$$

נסביר את המעבר הראשון: ניצר שני מקרים: בראשון  $b = 0$  ובשני  $b = 1$ .  
 $b_i = 0$ .

נמצא את התוחלת למקרה הזה:

$$\begin{aligned} \mathbf{E}[b_i \neq b] &= \int_x 2 \cdot p(x) \cdot (1 - p(x)) \cdot D(x) dx \\ &= \int_x 2 \cdot r(x) \cdot (1 - r(x)) \cdot D(x) dx \\ &= 2 \cdot \int_x r(x) \cdot D(x) dx - 2 \cdot \int_x r^2(x) \cdot D(x) dx \\ &\leq 2 \cdot R^* - 2 \cdot (R^*)^2 = 2 \cdot R^* \cdot (1 - R^*) \end{aligned}$$

7.3.2.2 המקרה הכללי

נניח:

- השכן הקרוב ביותר מתכנס ל- $x$  כאשר  $m \rightarrow \infty$ , ואז:

$$\mathbf{E}_x \mathbf{E}_{x_1, \dots, x_m} \min_i \|x_i - x\| \xrightarrow{m \rightarrow \infty} 0$$

- הסיווג בסביבה של  $x$  דומה לזה של  $x$ :

$$|p(x) - p(z)| \leq \alpha \cdot \|x - z\|$$

נוסיף עוד שתי הנחות מפשטות:

- $D(x)$  לא אפס על כל  $X$ , ואין נקודת מסה.
- $p(x)$  פונקציה רציפה.

שתי ההנחות האחרונות אינן הכרחיות, אך מפשטות מאוד את המודל. בהינתן  $x_1, \dots, x_m$ , נגדיר את  $\text{NN}_1^m(x)$  להיות השכן הקרוב ביותר ל- $x$ :

$$\text{NN}_1^m(x) = \arg \min_{x_i} \|x_i - x\|$$

**משפט** עבור כל  $x$ :

$$\text{NN}_1^m(x) \xrightarrow{m \rightarrow \infty} x$$

**הוכחה** לכל  $\varepsilon > 0$ , נגדיר:  $B_\varepsilon(x)$  - כדור ברדיוס  $\varepsilon$  שמרכזו  $x$ . אז:

$$\Pr[x_i \in B_\varepsilon(x)] = q > 0$$

אם נדגום  $m$  פעמים:

$$\Pr[x_1, \dots, x_m \notin B_\varepsilon(x)] = (1 - q)^m \xrightarrow{m \rightarrow \infty} 0$$

□

וזה נכון לכל  $\varepsilon > 0$ .

למעשה, הוכחנו כי:

$$\forall \varepsilon > 0, \delta > 0. \exists m_D(\varepsilon, \delta). \Pr[\text{NN}_1^m(x) \notin B_\varepsilon(x)] \leq \delta$$

**משפט**

$$r(\text{NN}_1^m(x)) \xrightarrow{m \rightarrow \infty} r(x)$$



□ הוכחה נובעת מרציפות  $p(x)$ . נגדיר:  $e_x^m$  - המאורע ש- $\text{NN}_1^m(x)$  עשה שגיאה עם מדגם בגודל  $m$  על הנקודה  $x$ . אזי:

$$\Pr [e_x^m] \xrightarrow{m \rightarrow \infty} r(x) \cdot (1 - r(x))$$

לכן:

$$\begin{aligned} \lim_{m \rightarrow \infty} \int \Pr [e_x^m] \cdot D(x) dx &= \int \lim_{m \rightarrow \infty} \Pr [e_x^m] \cdot D(x) dx \\ &= \int r(x) \cdot (1 - r(x)) \cdot D(x) dx \end{aligned}$$

## 7.4 שיטת $k$ שכנים קרובים ( $k$ -NN)

נגדיר:

$$\text{NN}_k^m = \left\{ \begin{array}{l} k \text{ nearest} \\ \text{neighbors of } m \end{array} \right\}$$

אזי  $\hat{p}(x) = \frac{\ell}{k}$  אם  $\ell$  מתוך  $k$  השכנים הקרובים חיוביים. נבחר:

1.

$$\frac{k}{m} \rightarrow 0$$

2.

$$k \rightarrow \infty$$

למשל:  $k = \sqrt{m}$  ו- $k = \log m$ .  
נרצה להראות:

1. כל  $k$  השכנים הקרובים ביותר מתכנסים ל- $x$ :

$$\forall z \in \text{NN}_k^m(x). z \rightarrow x$$

2.

$$\hat{p}(x) \rightarrow p(x)$$

### 7.4.1 מקרה פשוט

כל  $k$  הקרובים ביותר זהים ל- $x$  (כל השאר שונים מ- $x$ ).  
יש לנו  $k$  מטבעות עם הסתברות  $\hat{p}(x)$ :

$$\Pr [|\hat{p}(x) - p(x)| \geq \lambda] \leq 2 \cdot e^{-\lambda^2 \cdot k} \cdot \hat{p}(x)$$

לכל  $\lambda > 0$ , ומכיוון ש- $k \rightarrow \infty$ , מתקיים  $\hat{p}(x) \rightarrow p(x)$ .

**7.4.2 מקרה כללי**התכנסות ל- $x$ .כמו קודם, נגדיר את  $B_\varepsilon(x)$ . אזי:

$$\Pr[x_i \in B_\varepsilon(x)] = q > 0$$

מספר הנקודות  $z$  בתוחלת שיפלו בתוך  $B_\varepsilon(x)$  הוא  $q \cdot m$ . לכן:

$$\begin{aligned} \Pr[\text{at most } k-1 \text{ } x_i\text{'s in } B_\varepsilon(x)] &= \Pr\left[\left|\frac{z}{m} - q\right| \geq q - \frac{k-1}{m}\right] \\ &\leq 2 \cdot e^{-\left(q - \frac{k-1}{m}\right)^2 \cdot m} \xrightarrow{m \rightarrow \infty} 0 \end{aligned}$$

**7.5 מדידת המרחק**

המרחק תלוי בגודל של המימדים, או ביחידות שלהם (למשל ק"ג או גרמים, שניות או דקות). רצוי להביא את המדידות לסקלה אחידה. איך עושים את זה? למשל:

- ניח שהמדידות התקבלו בהתפלגות נורמלית:

$$\frac{v_i - \mu_i}{\sigma_i}$$

- ניח שהמדידות התקבלו בהתפלגות אחידה:

$$\frac{v_i}{\max_i - \min_i}$$

אנחנו צריכים לבחור מרחק, למשל:

$$L_2(x, y) = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$$

או לפי כל נורמה אחרת שנבחר.

**7.6 Locality Sensitive Hashing**

השאלה שתעניין אותנו: בהינתן  $x$  ו- $R$ , צריך למצוא  $y \in S$  כך ש- $\|x - y\| \leq R$  או להגיד שאין כזה, כלומר:

$$\forall y \in S. \|x - y\| > R$$

הגדרה תהי משפחה  $H$ . אם לכל  $p, q \in \mathbb{R}^d$ :

1. אם  $\|p - q\| \leq R$  אזי  $\Pr_H [h(p) = h(q)] \geq p_1$ .

2. אם  $\|p - q\| > R$  אזי  $\Pr_H [h(p) = h(q)] \leq p_2$ .

3.  $p_1 < p_2$ .

אזי נאמר כי  $H$  היא משפחה טובה. זו ההגדרה ל-*Locality Sensitive Hashing*.

דוגמה  $X = \{0, 1\}^d$ . נגדיר:

$$h_i(x_1, \dots, x_d) = x_i$$

אזי:

$$\Pr_H [h(x) = h(y)] = 1 - \frac{d_H(x, y)}{d}$$

כאשר  $d_H$  הוא מספר המינגטון:

$$d_H(x, y) = \# \text{ corr. diff}$$

### 7.6.1 שלב א' - Amplification

נגדיר:

$$g(x) = \langle h_1(x), \dots, h_k(x) \rangle$$

כאשר  $h_i$  נבחרים באקראי מתוך  $H$ .

עבור  $g$ :

1. אם  $\|p - q\| \leq R$ , אזי:

$$\Pr [g(p) = g(q)] \geq p_1^k$$

2. אם  $\|p - q\| > R$ , אזי:

$$\Pr [g(p) = g(q)] \leq p_2^k$$

נבחר את  $k$  כך ש- $p_2^k = \frac{1}{n}$ . אזי:

$$k = \frac{\ln n}{\ln 1/p_2}$$

לכן:

$$\begin{aligned} p_1^k &= p_1^{\frac{\ln n}{\ln 1/p_2}} \\ &= n^{-\frac{\ln 1/p_1}{\ln 1/p_2}} \\ &= n^{-\rho} \end{aligned}$$

כאשר:

$$\rho = \frac{\ln 1/p_1}{\ln 1/p_2}$$

**7.6.2 שלב ב'**

נבחר  $g_1, \dots, g_L$ , ונטען שההסתברות שקיים  $i$  כך ש- $g_i(p) = g_i(q)$  כאשר  $\|p - q\| \leq R$  גדולה.

ההסתברות עבור  $g$  יחיד היא  $1 - n^{-\rho}$ . לכן, ההסתברות עבור  $g_1, \dots, g_L$  היא:

$$(1 - n^{-\rho})^L \approx e^{-n^\rho \cdot L} = \delta$$

לכן:

$$L = n^\rho \cdot \ln \frac{1}{\delta}$$

**7.6.3 האלגוריתם**

בוחרים  $L$  פונקציות  $g_1, \dots, g_L$  כך ש- $g_i(x) = \langle h_{i,1}(x), \dots, h_{i,k}(x) \rangle$  כאשר  $h_{i,j} \in H$  אקראיים.

עבור כך  $g_i$ : לכל  $x$  מחשבים את  $g_i(x)$ , ומכניסים לטבלת Hash. יש לנו  $L$  טבלאות. בסך הכל, נתפוס  $O(L \cdot n)$  מקום.

חיפוש: בהינתן  $q$ , לכל  $i$  נחשב את  $g_i(q)$ , ונבדוק אם יש  $g_i(p) = g_i(q)$  בטבלה ה- $i$ . זמן החישוב יהיה:

$$\underbrace{\frac{1}{n} \cdot n \cdot L}_{\text{Bad points}} + \underbrace{L \cdot \left( \begin{array}{c} \text{Hash} \\ \text{lookup} \end{array} \right) \cdot \left( \begin{array}{c} \# \\ \text{Nearest} \\ \text{points} \end{array} \right)}_{\text{Good points}}$$

איך נעשה רדוקציות כאלה? למשל אם  $X = \{0, \dots, S-1\}^d$ . אזי כל  $x \in X$  יהיה בייצוג אונרי<sup>2</sup>.

במקרה הזה,  $L_1$  יהיה  $d_H$  (המינגטון) על  $\{0, 1\}^{d \cdot \log S}$ , כאשר:

$$L_1(x, y) = \sum_{i=1}^d (x_i - y_i)$$

**דוגמה** מרחק  $L_1$  מעל  $[0, 1]^d$ . נניח כי  $R \ll 1$ . נבחר באקראי  $s_1, \dots, s_d \in [0, 1]$  אזי:

$$h_{s_1, \dots, s_d}(x) = \langle \text{sign}(x_1 - s_1), \dots, \text{sign}(x_d - s_d) \rangle$$

לכן:

$$\Pr_s [\text{sign}(x_i - s_i) \neq \text{sign}(y_i - s_i)] = |x - y|$$

כאשר  $x, y \in [0, 1]^d$

<sup>2</sup>למשל, 3 ייוצג כך:  $\underbrace{111}_{3 \text{ times}} \underbrace{00 \dots 0}_{3 \text{ times}}$

:זאן

$$\begin{aligned}\Pr[h(x) = h(y)] &= \prod_{i=1}^d (1 - |x_i - y_i|) \\ &= 1 - \prod_{i=1}^d |x_i - y_i| \\ &= 1 - L_1(x, y)\end{aligned}$$

## פרק 8

# $^1$ VC Dimension

### 8.1 מודל PAC (חזרה)

במודל ה-PAC נתונה לנו  $D$ , ההתפלגות מעל הקלטים  $X$ . ההנחות שלנו לגביה:

1.  $D$  אינה משתנה (עם התקדמות הלמידה).

2.  $D$  אינה ידועה (ללומד).

3. הדוגמאות נדגמות מ- $D$  באופן i.i.d, גם עבור הלמידה וגם עבור הבדיקה.

יש לנו פונקציית מטרה  $c^* \in C$ . הדוגמאות הן מהצורה  $\langle x, c^*(x) \rangle$  כאשר  $x \in X$ . המטרה: למצוא  $h \in H$  כך שהשגיאה, המוגדרת על ידי:

$$\text{error}(h) = D[c^*(x) \neq h(x)]$$

תהיה קטנה ככל הניתן.

יש לנו שני פרמטרים:  $\varepsilon$  - פרמטר הדיוק, ו- $\delta$  - פרמטר הבטחון. מאורע של הצלחה מוגדר כך:

$$\text{error}(h) \leq \varepsilon$$

פרמטר הבטחון נותן לנו:

$$\Pr[\text{success}] \geq 1 - \delta$$

אלגוריתם PAC מקבל את  $\varepsilon$  ו- $\delta$ , ולפעמים יש לו גישה לאורקל שדוגם את  $X$  לפי  $D$ :

$$\mathbf{EX}(D, c^*) \mapsto \langle x, c^*(x) \rangle$$

אלגוריתם PAC  $A$  לומד משפחה  $C$  (קבוצת פונקציות מטרה) אם לכל  $c^* \in C$  ולכל התפלגות  $D$  (מעל  $X$ ), מוצא השערה  $h \in H$  כך ש:

$$\Pr[\text{error}(h) \leq \varepsilon] \geq 1 - \delta$$

נאמר כי האלגוריתם  $A$  יעיל אם הוא רץ בזמן פולינומי ב- $1/\varepsilon, 1/\delta$  וקידוד  $c^*$ .

<sup>1</sup>שיעור שהתקיים בתאריך 16.12.2012.

## 8.2 מימד VC

## 8.2.1 מוטיבציה

נתעניין במקרה שבו  $C = H$ . אנחנו נרצה לראות כמה דוגמאות  $m$  נחוצות כדי שהשערה עקבית (שגיאה 0 על המדגם) תלמד PAC. אם  $C$  מחלקה סופית, ראינו כי:

$$m \geq \frac{1}{\varepsilon} \cdot \ln \frac{|C|}{\delta}$$

מה קורה אם  $C$  אינסופית? החסם הזה לא אומר לנו כלום. האם זה אומר שאנחנו לא נצליח ללמוד מחלקות אינסופיות? ראינו כבר את הדוגמה של למידת מלבנים מקבילים לצירים, שבה:

$$m = O\left(\frac{1}{\varepsilon} \ln \frac{1}{\delta}\right)$$

אבל בהוכחה הזו הסתמכנו על תכונות של המלבנים. היינו רוצים למצוא דרך גנרית יותר כדי להבין מהו  $m$ .

## 8.2.2 הגדרות

נניח כי  $C$  מחלקת השערות מעל  $X$ . ניתן לחשוב על  $c \in C$  בשתי דרכים:

1.  $c(x) \in \{0, 1\}$ .

2.  $\{x \mid c(x) = 1\}$ .

עבור  $C$  מעל  $X$  ו- $S \subseteq X$  סופית, נגדיר:

$$\Pi_C(S) = \{c \cap S \mid c \in C\}$$

או באופן שקול, אם  $S = \{x_1, \dots, x_m\}$ :

$$\Pi_C(S) = \{\langle c(x_1), \dots, c(x_m) \rangle \mid c \in C\}$$

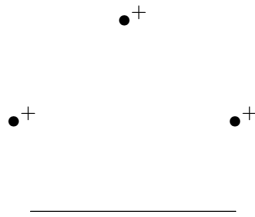
אנחנו נתעניין בגודל שצ  $\Pi_C(S)$ . נקרא ל- $\Pi_C(S)$  הטלה של  $C$  על קבוצה  $S$ . נתעניין בשאלה הבאה: האם  $C$  משרה את כל הפונקציות מעל  $S$ ? כלומר, האם  $|\Pi_C(S)| = 2^m$ ?

**הגדרה** מחלקה  $C$  פותצת את  $S$  אם  $|\Pi_C(S)| = 2^m$  כאשר  $m = |S|$ .

**הגדרה** נגדיר את מימד VC:

$$\text{VCdim}(C) = \max \{d \mid \exists S. |S| = d \wedge |\Pi_C(S)| = 2^d\}$$

אם התנאי מתקיים לכל  $d$ , אזי  $\text{VCdim}(C) = \infty$ .



איור 8.1: סיווג אפשרי של 3 נקודות שאינן על אותו הישר, עם המפריד הלינארי שלהן

### 8.2.3 דוגמאות

#### 8.2.3.1 סיפא על קו

נגדיר את  $C_1$  להיות מחלקת הסיפאות על הקווים,  $X = [0, 1]$ , ונגדיר:

$$c_\alpha(x) = \begin{cases} 0 & x < \alpha \\ 1 & x \geq \alpha \end{cases}$$

נשים לב כי  $C_1$  אינסופית (אפילו לא בת-מניה). נראה כי  $\text{VCdim}(C_1) = 1$ .

**שלב ראשון** נוכיח כי  $\text{VCdim}(C_1) \geq 1$ .  
נבחר  $S = \{1/2\}$ . אזי:

$$c_{3/4}\left(\frac{1}{2}\right) = 0$$

$$c_{1/4}\left(\frac{1}{2}\right) = 1$$

לכן,  $|\Pi_C(\{1/2\})| = 2$ . לכן,  $\text{VCdim}(C_1) \leq 1$ .

**שלב שני** נראה כי  $\text{VCdim}(C_1) < 2$ .

יהיו  $x, y \in [0, 1]$ . נניח כי  $y \geq x$ . אזי לא קיימת  $c \in C_1$  כך ש- $c(x) = 1$  ו- $c(y) = 0$ .  
לכן,  $\text{VCdim}(C_1) < 2$ .  
לכן,  $\text{VCdim}(C_1) = 1$ .

#### 8.2.3.2 מפריד לינארי במישור

נגדיר את  $C_2$  להיות מחלקת המפרידים הלינארים במישור,  $X = \mathbb{R}^2$ .  
אם  $(x_1, x_2) = x \in X$ , עבור  $w = (\alpha_1, \alpha_2, \theta)$  נגדיר:

$$c_w(x) = 1 \Leftrightarrow \sum_{i=1}^2 \alpha_i \cdot x_i \geq \theta$$

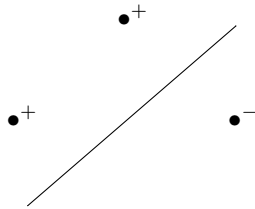
נראה כי  $\text{VCdim}(C_2) = 3$ .

1. נראה כי  $\text{VCdim}(C_2) \geq 3$ .

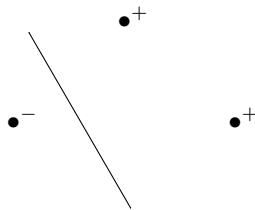
נבחר שלוש נקודות שאינן על אותו הישר. כפי שניתן לראות באיורים 8.1, 8.2, 8.3, 8.4, 8.5, 8.6, 8.7 ו-8.8, לכל סיווג אפשרי של שלוש הנקודות נוכל למצוא מפריד לינארי.

לכן,  $\text{VCdim}(C_2) \geq 3$ .

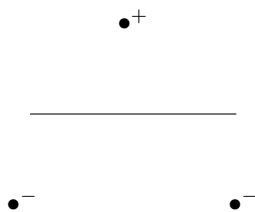




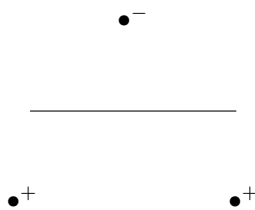
איור 8.2: סיווג אפשרי של 3 נקודות שאינן על אותו הישר, עם המפריד הלינארי שלהן



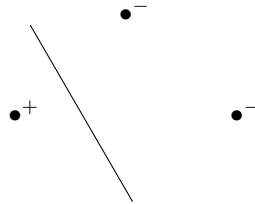
איור 8.3: סיווג אפשרי של 3 נקודות שאינן על אותו הישר, עם המפריד הלינארי שלהן



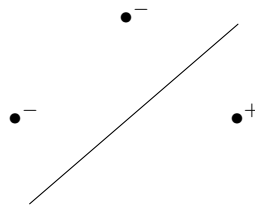
איור 8.4: סיווג אפשרי של 3 נקודות שאינן על אותו הישר, עם המפריד הלינארי שלהן



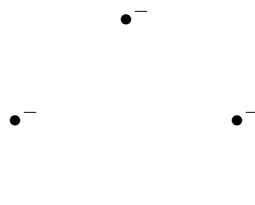
איור 8.5: סיווג אפשרי של 3 נקודות שאינן על אותו הישר, עם המפריד הלינארי שלהן



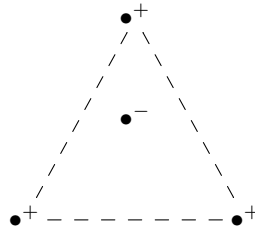
איור 8.6: סיווג אפשרי של 3 נקודות שאינן על אותו הישר, עם המפריד הלינארי שלהן



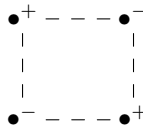
איור 8.7: סיווג אפשרי של 3 נקודות שאינן על אותו הישר, עם המפריד הלינארי שלהן



איור 8.8: סיווג אפשרי של 3 נקודות שאינן על אותו הישר, עם המפריד הלינארי שלהן



איור 8.9: סיווג בלתי ניתן לסיפוק של 4 נקודות במישור



איור 8.10: סיווג בלתי ניתן לסיפוק של 4 נקודות במישור

2. נראה כי  $VCdim(C_2) < 4$ .

בהינתן קבוצה כלשהי של 4 נקודות במישור, אנחנו נמצאים באחד המקרים שמתוארים באיורים 8.9, 8.10 או 8.11. כל אחד מהאיורים מראה סיווג בלתי ניתן לסיפוק על ידי  $C_2$ .

לכן,  $VCdim(C_2) < 4$ .

לכן,  $VCdim(C_2) = 3$ .

נוכיח בהמשך כי לכל מימד  $d \geq 2$ , מפריד לינארי במימד  $d$  מקיים  $VCdim(C^d) = d + 1$ .

### 8.2.3.3 מלבנים מקבילים לצירים

נגדיר את  $C_3$  להיות המלבנים המקבילים לצירים. נראה כי  $VCdim(C_3) = 4$ .

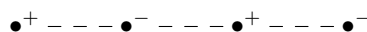
**הערה** לא ניתן לנתח כל קבוצה בגודל 4. למשל, הקבוצה שמתוארת באיור 8.12 אינה ניתנת לניתוח.

אבל אנחנו יכולים לבחור 4 נקודות כלשהן, ולא בהכרח את הנקודות האלה. נבחר 4 נקודות כמו באיור 8.13. קל לראות שלכל אחת מההצבות לסיווגי הנקודות ניתן לייצר מלבן חוסם. לכן,  $VCdim(C_3) \geq 4$ .

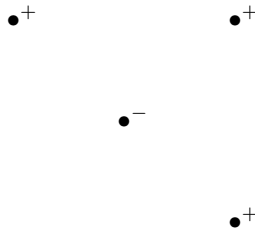
עבור 5 נקודות: לכל 5 נקודות, קיימת אחת שאינה קיצון באף מימד. נסמן אותה כשליית, ואת שאר הנקודות כחיוביות. לא קיימת  $c \in C_3$  שתנתח את 5 הנקודות. לכן,  $VCdim(C_3) < 5$ .

לכן,  $VCdim(C_3) = 4$ .

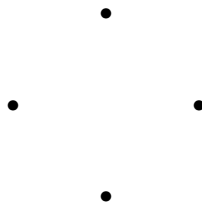
<sup>2</sup>ראה קטע 8.2.5.3.



איור 8.11: סיווג בלתי ניתן לסיפוק של 4 נקודות במישור



איור 8.12: דוגמה למקרה שלא ניתן לנתח קבוצה בגודל 4



איור 8.13: דוגמה לנקודות שכל הצבה שלהן מנתצת את  $C_3$

**8.2.3.4 מספר סופי של אינטרוולים**

נגדיר את  $C_4$  להיות מספר סופי של אינטרוולים. כלומר: בהינתן  $c \in C_4$ , קיימת סדרת קטעים  $\{[a_i, b_i]\}_{i=1}^{k_c}$  כך ש:

$$c(x) = 1 \Leftrightarrow x \in \bigcup_{i=1}^{k_c} [a_i, b_i]$$

נראה ש- $\text{VCdim}(C_4) = \infty$ . לכל  $d$  נקודות ולכל הצבה אפשר לשים אינטרוול קטן סביב הנקודות החיוביות. כלומר, לכל הצבה יש פונקציה ב- $C_4$  שהיא עקבית איתה. לכן,  $\text{VCdim}(C_4) = \infty$ .

**8.2.3.5 פוליגון קונבקסי במישור**

נגדיר את  $C_5$  להיות פוליגון קונבקסי במישור. נראה ש- $\text{VCdim}(C_5) = \infty$ . ניקח  $d$  נקודות על מעגל היחידה. לכל הצבה בנקודות, נוכל לחבר את הנקודות החיוביות ולקבל פוליגון קונבקסי. כל הנקודות השליליות יהיו מחוץ לפוליגון. לכן, לכל  $d$  נקודות ולכל הצבה יש פוליגון קונבקסי שעקבי עם הצבה. לכן,  $\text{VCdim}(C_5) = \infty$ .

**8.2.4 חסם תחתון על גודל הדגימה**

**משפט** אם  $\text{VCdim}(C) = d + 1$  אזי:

$$m\left(\varepsilon, \delta = \frac{1}{2}, d + 1\right) \geq \frac{d}{16 \cdot \varepsilon} = \Omega\left(\frac{d}{\varepsilon}\right)$$

**הוכחה** נבחר  $T = \{z_0, \dots, z_d\}$  כך ש- $|T| = 2^{d+1}$ . נבנה את ההתפלגות  $D$ :

$$D(x) = \begin{cases} 1 - 8 \cdot \varepsilon & x = z_0 \\ \frac{8 \cdot \varepsilon}{d} & x = z_i, i \geq 1 \\ 0 & \text{otherwise} \end{cases}$$

נבחר פונקציית מטרה:

$$c^*(x) = \begin{cases} 0 \text{ or } 1 \text{ with probability } 1/2 & x \in T \\ \text{don't care} & x \notin T \end{cases}$$

קיימת  $c \in C$  שזהה ל- $c^*$  ב- $T$ : נשתמש ב- $c$  עבור  $x \notin T$ . נשים לב כי  $c^*$  אינה פונקצייה אקראית, אלא הבנייה שלה אקראית. נגדיר את הקבוצה  $SEEN$  כקבוצת כל הנקודות ב- $T$  שראינו (במהלך הלמידה). כמו כן, נגדיר:  $UNSEEN = T \setminus SEEN$ . לכל  $z \in UNSEEN$ , הסתברות השגיאה היא בדיוק  $1/2$ . בהנחה שראינו את  $z_0$ :

$$\begin{aligned} D[\text{error}] &\geq \frac{1}{2} \cdot |UNSEEN| \cdot \frac{8 \cdot \varepsilon}{d} \\ &= \frac{4 \cdot \varepsilon}{d} \cdot |UNSEEN| \end{aligned}$$

כמו כן:

$$\mathbf{E}[|SEEN|] \leq m \cdot 8 \cdot \varepsilon$$

נניח בשלילה כי  $m < \frac{d}{16 \cdot \varepsilon}$ . אזי  $\mathbf{E}[|SEEN|] < \frac{d}{2}$ . לכן:

$$\Pr \left[ |SEEN| \leq \frac{d}{2} \right] \geq \frac{1}{2}$$

לכן, בהסתברות  $1/2$ ,  $|UNSEEN| \geq \frac{d}{2}$ , ולכן  $D[\text{error}] \geq 2 \cdot \varepsilon$ . בסתירה להגדרת  $m$ .  $\square$

### 8.2.5 עוד דוגמאות

#### 8.2.5.1 פונקציית Parity

נגדיר את  $C_6$  להיות פונקציות Parity. כלומר,  $X = \{0, 1\}^n$  ועבור  $S \subseteq \{1, \dots, n\}$  מגדירים:

$$\chi_S(x) = \bigoplus_{i \in S} x_i$$

נראה ש- $\text{VCdim}(C_6) = n$ .

1. נראה ש- $\text{VCdim}(C_6) \geq n$ .

נבחר את  $n$  וקטורי הבסיס  $\{e_i\}_{i=1}^n$ , כאשר:

$$e_i = \left( \underbrace{0, \dots, 0}_{i-1}, 1, \underbrace{0, \dots, 0}_{n-i} \right)$$

לכל הצבה  $b_1, \dots, b_n$  עבור  $e_1, \dots, e_n$  נבחר  $S = \{i \mid b_i = 1\}$ , ואז:

$$\chi_S(e_j) = \begin{cases} 1 & j \in S \\ 0 & j \notin S \end{cases}$$

כאשר  $j \in S \Leftrightarrow b_j = 1$ .

לכן,  $\text{VCdim}(C_6) \geq n$ .

2. נראה ש- $\text{VCdim}(C_6) < n + 1$ .

יש שתי דרכים להראות את זה:

(א) משיקולים קומבינטוריים:  $|C_6| = 2^n$ .

כמו כן, באופן כללי,  $\text{VCdim}(C) \leq \log_2 |C|$ . למה? אם ל- $C$  יש  $\text{VCdim}(C) = d$ , אז יש לפחות  $2^d$  פונקציות שונות.

(ב) בהינתן וקטורים  $x_1, \dots, x_{n+1}$ , קיים וקטור  $x_i$  כך ש- $x_i = \bigoplus_{j \in T} x_j$ . אם נבחר  $b_j = 0$  לכל  $j \in T$ , אז בהכרח  $\chi_S(x_i) = 0$ . לכן, ההצבה  $b_i = 1$  אינה אפשרית.

לכן,  $\text{VCdim}(C_6) < n + 1$ .

לכן,  $\text{VCdim}(C_6) = n$ .

### 8.2.5.2 OR של ליטרלים

נגדיר את  $C_7$  להיות OR של ליטרלים כאשר  $X = \{0, 1\}^n$ . נראה ש- $\text{VCdim}(C_7) = n$ .

1. נראה ש- $\text{VCdim}(C_7) \geq n$ .

נבחר את וקטורי הבסיס  $e_1, \dots, e_n$ . בהינתן הצבה  $b_1, \dots, b_n$ , נגדיר את  $S$  כמו קודם:  $S = \{i \mid b_i = 1\}$ . ואז:

$$c_S(x) = \bigvee_{i \in S} x_i$$

2. נראה ש- $\text{VCdim}(C_7) \leq n$ .

לצורך ההוכחה, נשתמש בטענת עזר:

**טענה** אם למחלקה  $C$  יש אלגוריתם Online שעושה לכל היתר  $d$  שגיאות, אזי  $\text{VCdim}(C) \leq d$ .

**הוכחה** נניח בשלילה כי  $\text{VCdim}(C) \geq d + 1$ . אזי קיימות  $d + 1$  נקודות  $S$  שמקבלות את כל ההצבות. נסמן:  $S = \{x_1, \dots, x_{d+1}\}$ .

נרץ את האלגוריתם על  $S$ , ותמיד נגיד לו שטעה. האלגוריתם יעשה  $d + 1$  שגיאות, בסתירה להנחה.  $\square$

**למה** בהינתן 3 וקטורים (או יותר), ישנם 2 וקטורים במרחק 2 לפחות.

**הוכחה** יהיו  $z_1, z_2, z_3 \in \{0, 1\}^n$  3 וקטורים שונים. נניח כי  $z_1$  ו- $z_2$  נבדלים בבית  $i$ , ונניח כי  $z_2$  ו- $z_3$  נבדלים בבית  $j$ .

אם  $i = j$ , אזי  $z_1 = z_3$ , בסתירה להנחה כי הווקטורים שונים. לכן,  $z_1$  ו- $z_3$  נבדלים בביתים  $i$  ו- $j$ .  $\square$

נראה כעת כי  $\text{VCdim}(C_7) \leq n$ : נראה שיש סדר כך שהאלגוריתם  $^3\text{ELIM}$  מבצע לכל היותר  $n$  שגיאות.

בהינתן  $n+1$  נקודות, נבחר את  $a$  ו- $b$  להיות שני וקטורים במרחק 2 לפחות. נרץ את  $\text{ELIM}$  על הסדרה הבאה:  $a, b$  ואז כל האחרים. כאשר עוברים על  $a$ , מקבלים שגיאה ופוסלים לכל היותר  $n$  ליטרלים. כשעוברים על  $b$  פוסלים עוד 2 ליטרלים. כשעוברים על הווקטורים האחרים, עושים עוד  $n-2$  איטרציות, שבכל אחת פוסלים לכל היותר ליטרל בודד.

לכן, קיבלנו לכל היותר  $n$  שגיאות. לכן,  $\text{VCdim}(C_7) \leq n$ .

לכן,  $\text{VCdim}(C_7) = n$ .

### 8.2.5.3 מפריד לינארי במימד $n$

נגדיר את  $C_8$  להיות מפריד לינארי במימד  $n$ ,  $X = \mathbb{R}^n$ ,  $(n \geq 2)$ . עבור  $w = (\alpha_1, \dots, \alpha_n, \theta) \in C_8$  נגדיר את  $c_w$ :

$$c_w(x) = 1 \Leftrightarrow \sum_{i=1}^n \alpha_i \cdot x_i \geq \theta$$

נוכיח כי  $\text{VCdim}(C_8) = n+1$ .

1. נראה כי  $\text{VCdim}(C_8) \geq n+1$ .

נגדיר:  $E = \{\vec{0}, \vec{e}_1, \dots, \vec{e}_n\}$ . בהינתן  $b_0, \dots, b_n$  נבחר:

$$\alpha_i = \begin{cases} +1 & b_i = 1 \\ -1 & b_i = 0 \end{cases}$$

$$\theta = \begin{cases} -\frac{1}{2} & b_0 = 1 \\ +\frac{1}{2} & b_0 = 0 \end{cases}$$

נסמן:  $w = (\alpha_1, \dots, \alpha_n, \theta)$ . אזי:

$$c_w(e_i) = 1 \Leftrightarrow \alpha_i \geq \theta \Leftrightarrow b_i = 1$$

$$c_w(\vec{0}) = 1 \Leftrightarrow 0 \geq \theta \Leftrightarrow b_0 = 1$$

$\square$

לכן,  $c_w$  עיקבית עם  $b_0, \dots, b_n$ .

2. נראה כי  $\text{VCdim}(C_8) < n+2$ .

<sup>3</sup>ראה אלגוריתם 1.3.

**הגדרה** תת-קבוצה  $A \subseteq \mathbb{R}^d$  היא קונבקסית אם:

$$\forall x_1, x_2 \in A. \forall \lambda \in [0, 1]. x_1 + (1 - \lambda)x_2 \in A$$

ה- $Convex Hull$  של קבוצת נקודות  $S$  הוא הקבוצה הקונבקסית הקטנה ביותר שמכילה את  $S$ . ה- $Convex Hull$  של  $S$  מסומן ב- $Conv(S)$ .

**משפט Radon** תהי  $E \subseteq \mathbb{R}^d$  כך ש- $|E| = d + 2$ . אזי קיימת  $S \subseteq E$  כך ש- $S \neq \emptyset$  וגם  $E \setminus S \neq \emptyset$  וגם  $Conv(S) \cap Conv(E \setminus S) \neq \emptyset$ .

**הוכחה נסמן:**  $E = \{x_0, \dots, x_{d+1}\}$  אי קיימים  $\alpha_0, \dots, \alpha_{d-1}$  כך ש- $\sum_{i=0}^{d-1} \alpha_i = 1$  ו- $x_i = \vec{0}$  מכיוון ש- $\mathbb{R}^d$  ממימד  $d$  ו- $|E| = d + 2$ , אזי אפשר להוסיף את האילוץ כי  $\sum_{i=0}^{d+1} \alpha_i = 0$ .

קיבלנו בסך הכל  $d + 1$  משוואות עבור  $d + 1$  נעלמים  $(\alpha_0, \dots, \alpha_{d-1})$ , ולכן קיים פתרון שאינו טריוויאלי.

נניח בה"כ כי  $\alpha_0, \dots, \alpha_p$  חיוביים ו- $\alpha_{p+1}, \dots, \alpha_{d+1}$  שליליים. נסמן:

$$\alpha = \sum_{i=0}^p \alpha_i > 0$$

עבור  $i \leq p$  נגדיר:

$$\beta_i = \frac{\alpha_i}{\alpha} > 0$$

עבור  $i \geq p + 1$  נגדיר:

$$\gamma_i = \frac{-\alpha_i}{\alpha} > 0$$

נסמן:

$$y = \sum_{i=0}^p \beta_i \cdot \vec{x}_i = \sum_{i=p+1}^{d+1} \gamma_i \cdot \vec{x}_i$$

נסיק כי:

$$\sum_{i=0}^p \beta_i = 1 = \sum_{i=p+1}^{d+1} \gamma_i$$

נקבע:  $S = \{x_0, \dots, x_p\}$ ,  $E \setminus S = \{x_{p+1}, \dots, x_{d+1}\}$  אזי  $y \in Conv(S) \cap Conv(E \setminus S) \neq \emptyset$ .  
 $\square$

נוכיח כעת כי  $VCdim(C_8) < n + 2$ : נניח בשלילה שלא. תהי  $E$  קבוצה של  $n + 2$  נקודות. לפי משפט Radon, קיימת  $S \subseteq E$  כך ש- $Conv(S) \cap Conv(E \setminus S) \neq \emptyset$ .

ניתן ל- $S$  סיווג חיובי ול- $E$  סיווג שלילי. אזי קיים  $c_w$  שמפריד בין  $S$  לבין  $E \setminus S$ . נסתכל על  $x \in Conv(S) \cap Conv(E \setminus S)$ :

•  $x$  חיובי כי  $S$  חיובית.

•  $x$  שלילי כי  $E \setminus S$  שלילית.

סתירה! לכן, לא קיימות  $n + 2$  נקודות שניתנות לניתוח. לכן,  $VCdim(C_8) < n + 2$ .

לכן,  $VCdim(C_8) = n + 1$ .



## פרק 9

# מימד VC (המשך)<sup>1</sup>

### 9.1 חזרה

במודל PAC, יש לנו התפלגות  $D$  על הקלטים  $X$ , ומחלקת השערות  $C$ . פונקציית המטרה שלנו היא  $c^* \in C$ , ומניחים כי  $C = H$ . האלגוריתם מוצא  $h \in H$  כך שבהסתברות  $1 - \delta$ ,  $\text{error}(h, c^*) \leq \epsilon$ , כאשר  $\delta$  ו- $\epsilon$  הם קלטים לאלגוריתם. עבור  $S \subseteq X$  ומחלקה  $C$  נגדיר:

$$\Pi_C(S) = \{c \cap S \mid c \in C\}$$

או באופן שקול:

$$\Pi_C(S) = \{(c(x_1), \dots, c(x_m)) \mid c \in C\}$$

כאשר  $S = \{x_1, \dots, x_m\}$ . אזי  $1 \leq |\Pi_C(S)| \leq 2^m$ . נרמא שמחלקה  $C$  פנתצת (*Shatters*) את  $S$  אם  $|\Pi_C(S)| = 2^m$ . נגדיר את עיפד  $VC$ :

$$\text{VCdim}(C) = \max \{d \mid \exists S. |S| = d \wedge |\Pi_C(S)| = 2^d\}$$

אם לכל  $d$  קיימת קבוצה  $S$  בגודל  $d$  ש- $C$  מנתצת, אזי:

$$\text{VCdim}(C) = \infty$$

אם  $d = \text{VCdim}(C)$ , עבור מחלקה סופית:

$$|C| \geq |\Pi_C(S)|$$

לכן:

$$|C| \geq 2^d$$

ולכן:

$$\log_2(|C|) \geq d$$

---

<sup>1</sup>סיכום לשיעור שהתקיים בתאריך 23.12.2012.

## 9.2 חסמים לגודל הזגימה

בשיעור הקודם ראינו חסם תחתון לגודל הזגימה:

$$m \left( \varepsilon, \delta = \frac{1}{2}, d \right) \geq \Omega \left( \frac{d}{\varepsilon} \right)$$

היום נראה חסם עליון.

נדגום את  $S$ . מספר ההשערות ב- $\Pi_C(S)$  הוא סופי. לכן, יש לנו מחלקת השערות סופי. לכן:

$$m \leq \frac{1}{\varepsilon} \cdot \log \frac{|\Pi_C(S)|}{\delta}$$

הבעיה שלנו בטענה הזאת היא שקיבלנו את  $S$ , ורק אחר כך קבענו את החסם. החסם לא אמור להיות תלוי במדגם שלנו.

**הגדרה** נגדיר את קבוצת ההשערות הרעות:

$$B_\varepsilon(c^*) = \{h \in H \mid \text{error}(h, c^*) > \varepsilon\}$$

נרצה להראות שאם גודל המדגם  $m$  מספיק גדול, אזי  $h \notin B_\varepsilon(c^*)$  בהסתברות  $1 - \delta$ .

**הגדרה** קבוצת נקודות  $S$  היא  $\varepsilon$ -Net עבור  $c^* \in C$  והתפלגות  $D$  אם:

$$\forall h \in B_\varepsilon(c^*). \exists x \in S. c^*(x) \neq h(x)$$

נראה שבהסתברות  $1 - \delta$ ,  $S$  היא  $\varepsilon$ -Net, ולכן זוהי למידת PAC.

נגדיר מדגם  $S = S_1 \uplus S_2$  כאשר  $|S_1| = m = |S_2|$ . נגדיר את המאורע  $A$ :  $S_1$  הוא לא  $\varepsilon$ -Net. נגדיר את המאורע  $B$ : ישנה פונקציה  $h \in B_\varepsilon(c^*)$  שהיא עקבית עם  $S_1$ , ול- $h$  יש לפחות  $\frac{\varepsilon \cdot m}{2}$  שגיאות על  $S_2$ . ברור ש- $A \subseteq B^c$ . כמו כן,  $\Pr[B \mid A] \geq 1/2$ . נחשב:

$$\begin{aligned} \Pr[B] &= \Pr[B \mid A] \cdot \Pr[A] \\ &\geq \frac{1}{2} \cdot \Pr[A] \\ 2 \cdot \Pr[B] &\geq \Pr[A] \end{aligned}$$

נחסום את  $\Pr[B]$ . נגדיר:

$$\begin{aligned} F &= \Pi_C(S_1 \cup S_2) \\ \text{ER}(h) &= \{x \in S_1 \cup S_2 \mid c^*(x) \neq h(x)\} \end{aligned}$$

המאורע  $A$  אומר  $\text{ER}(h) \cap S_1 = \emptyset$ . המאורע  $B$  אומר  $\text{ER}(h) \cap S_1 \neq \emptyset$ .  
 המאורע  $A$  אומר  $S_2 = \text{ER}(h)$  כאשר  $|\text{ER}(h)| \geq \frac{m-\varepsilon}{2}$ .  
 יש לנו בעיה קומבינטורית: יש לנו  $2 \cdot m$  כדורים, שמתוכם לפחות  $\frac{m-\varepsilon}{2}$  שחורים (השאר לבנים). נסמן ב- $\ell$  את מספר הכדורים השחורים.

מספר הקונפיגורציות לצביעת הכדורים:  $\binom{2 \cdot m}{\ell}$ . מספר הקונפיגורציות בהן  $\ell$  רק ב- $S_2$ :  $\binom{m}{\ell}$ . לכן, ההסתברות שכל  $\ell$  הכדורים יהיו ב- $S_2$ :

$$\begin{aligned} \frac{\binom{m}{\ell}}{\binom{2 \cdot m}{\ell}} &= \prod_{i=0}^{\ell-1} \frac{m-i}{2 \cdot m-i} \\ &= \prod_{i=0}^{\ell-1} \left( \frac{1}{2} - \frac{i/2}{2 \cdot m-i} \right) \leq \frac{1}{2^\ell} \end{aligned}$$

לכן:

$$\Pr[B] \leq |F| \cdot 2^{-\ell} \leq |F| \cdot 2^{-\frac{m-\varepsilon}{2}} \leq \delta$$

לכן:

$$m = O\left(\frac{1}{\varepsilon} \cdot \ln \frac{|F|}{\delta}\right)$$

כאשר

$$|F| = \left| \Pi_C \left( \underbrace{S_1 \cup S_2}_{2 \cdot \text{mpoints}} \right) \right|$$

השלב הבא יהיה לחסום את גודל  $F$  כפונקציה של מימד ה-VC. נגדיר:

$$J(m, d) = J(m-1, d) + J(m-1, d-1)$$

כאשר:

$$\forall d \quad J(0, d) = 1$$

$$\forall m \quad J(m, 0) = 1$$

לכן, הנוסחה הסגורה של  $J(m, d)$  היא:

$$J(m, d) = \sum_{i=0}^d \binom{m}{i}$$

**טענה** אם  $\text{VCdim}(C) = d$  ו- $|S| = m$ , אזי:

$$|\Pi_C(S)| \leq J(m, d)$$

**הוכחה** נוכיח באינדוקציה של  $m$  ו- $d$  (למעשה באינדוקציה על  $m+d$ ).  
 בסיס: אם  $m=1$  או  $d=1$  אנחנו בסדר (מהגדרת  $J$ ).  
 צעד האינדוקציה: נניח נכונות לכל  $m'$  ו- $d'$  כך ש- $m'+d' < m+d$ . כלומר,  $|\Pi_C(S)| \leq J(m', d')$  כאשר  $|S| = m'$  ו- $\text{VCdim}(C) = d'$ .

נוכח עבור  $m$  ו- $d$ : נסמן:

$$\begin{aligned} S &= \{x_1, \dots, x_m\} \\ C_S &= \Pi_C(S) \end{aligned}$$

נסתכל על  $T = \{x_1, \dots, x_{m-1}\}$ . אזי:

$$\Pi_C(T) = \{y \mid y_0 \in \Pi_C(S) \vee y_1 \in \Pi_C(S)\}$$

נסמן:

$$C_\star = \{y \in \Pi_C(T) \mid y_0 \in \Pi_C(S) \wedge y_1 \in \Pi_C(S)\}$$

נראה כי  $|C_S| = |C_T| + |C_\star|$ . מהאינדוקציה:

$$|C_T| \leq J(m-1, d)$$

נראה ש- $|C_\star| \leq J(m-1, d-1)$ . אם  $C_\star$  מנתצת את  $\{x_1, \dots, x_i\}$ , אזי  $C$  מנתצת את  $\{x_1, \dots, x_i, x_m\}$ , כי כל הצבה ל- $\{x_1, \dots, x_i\}$  ב- $C_\star$  יש לה שני המשכים ב- $C$  עבור  $x_m$ . מכיוון ש- $\text{VCdim}(C) = d$ ,  $\text{VCdim}(C_\star) \leq d-1$ . לכן, מהנחת האינדוקציה:

$$|C_\star| \leq J(m-1, d-1)$$

בסך הכל:

$$\begin{aligned} |C_S| &= |C_T| + |C_\star| \\ &\leq J(m-1, d) + J(m-1, d-1) \\ &= J(m, d) \end{aligned}$$

□

למעשה, הראנו כי  $|F| \leq J(2 \cdot m, d)$ . נסתכל על  $J(m, d)$ :

$$J(m, d) \leq \begin{cases} 2^m & d \geq m \\ 2 \cdot m^d & d < m \end{cases}$$

מכך ש- $m = O\left(\frac{1}{\varepsilon} \cdot \ln \frac{|F|}{\delta}\right)$  נסיק כי:

$$M = O\left(\frac{1}{\varepsilon} \cdot \ln \frac{1}{\delta} + \frac{d}{\varepsilon} \cdot \ln \frac{d}{\varepsilon}\right)$$

## 9.3 סיבוכיות Radamacker

### 9.3.1 ממוצעי Radamacker

**הגדרה** עבור  $S = \{x_1, \dots, x_m\}$  ומחלקה  $H$  כאשר  $H = \{h: X \rightarrow \{\pm 1\}\}$ , נגדיר את ממוצעי Radamacker:

$$R_S(H) = \mathbf{E}_\sigma \left[ \max_{h \in H} \frac{1}{m} \cdot \sum_{i=1}^m \sigma_i \cdot h(x_i) \right]$$

כאשר:

$$\forall i. \Pr[\sigma_i = 1] = \frac{1}{2} = \Pr[\sigma_i = -1]$$

באופן כללי, עבור  $h \in H$  כלשהו,  $\mathbf{E}[\sigma_i \cdot h(x_i)] = 0$ . במקרה שלנו,  $\sigma_i = \pm 1$  ו- $h(x_i) = \pm 1$  ולכן:

$$\sigma_i = h(x_i) \Leftrightarrow \sigma_i \cdot h(x_i) = 1$$

מה הקשר של זה למימד VC? אם  $H$  מנתצת את  $S$ :

$$\forall \sigma_i. \exists h \in H. \frac{1}{m} \cdot \sum_{i=1}^m \sigma_i \cdot h(x_i) = 1$$

לכן,  $R_S(H) \leq 1$ . אם  $H$  מכילה  $h$  וגם  $\bar{h}$ , אזי  $R_S(H) \geq 0$ .

**דוגמה**  $H_2 = \{h_0, h_2\}$ ,  $H_1 = \{h_0, h_1\}$  כאשר:

$$h_0 = -1, \dots, -1$$

$$h_1 = 1, \dots, 1$$

$$h_2 = -1, \dots, -1, 1$$

אזי:

$$R_S(H_1) \approx \frac{\sqrt{m}}{m}$$

$$R_S(H_2) = \frac{1}{m}$$

נשים לב כי  $\text{VCdim}(H_1) = 1 = \text{VCdim}(H_2)$ , אבל  $R_S(H_1) \neq R_S(H_2)$ .

**הגדרה**

$$R_D(H) = \mathbf{E}_{S \sim D} [R_S(H)]$$

### 9.3.2 אי שוויון McDiarmid

יהיו  $X_1, \dots, X_m$  מ"מ ב"ת. תהי פונקציה. נניח כי:

$$\forall i. |\Phi(x_1, \dots, x_i, \dots, x_m) - \Phi(x_1, \dots, x'_i, \dots, x_m)| \leq c_i$$

אזי לכל  $\varepsilon > 0$ :

$$\Pr[\Phi(x) > \mathbf{E}[\Phi(x)] + \varepsilon] \leq \exp\left(\frac{-2 \cdot \varepsilon^2}{\sum_{i=1}^m c_i^2}\right)$$

**9.3.3 סיבוכיות Radamacker**משפט בהסתברות  $1 - \delta$ , לכל  $h \in H$ :

$$\begin{aligned} \text{error}_D(h) &\leq \text{error}_S(h) + R_D(H) + \sqrt{\frac{\ln 2/\delta}{m}} \\ &\leq \text{error}_S(h) + R_S(D) + 3 \cdot \sqrt{\frac{\ln 2/\delta}{m}} \end{aligned}$$

כאשר  $|S| = 2 \cdot m$ .

הוכחה נגדיר:

$$\text{MAXGAP}(S) = \max_{h \in H} |\text{error}_D(h) - \text{error}_S(h)|$$

צעד 1: נראה שבהסתברות  $1 - \delta$ ,  $\text{MAXGAP}(S) \leq \varepsilon$ , במקרה שלנו:

$$c_i = \frac{1}{|S|} = \frac{1}{m}$$

לכן, בהסתברות  $1 - \delta$ :

$$\text{MAXGAP}(S) \leq \mathbf{E}[\text{MAXGAP}(S)] + \underbrace{\sqrt{\frac{\ln 2/\delta}{m}}}_{\varepsilon}$$

צעד 2: נראה כי  $R_D(S) \geq \mathbf{E}_S[\text{MAXGAP}(S)]$  נוסף דגימה  $S'$  בגודל  $m$ . אז:

$$\mathbf{E}[\text{error}_{S'}(h)] = \text{error}_D(h)$$

לכן:

$$\begin{aligned} \mathbf{E}[\text{MAXGAP}(S)] &= \mathbf{E}_S \left[ \max_{h \in H} \mathbf{E}_{S'}[\text{error}_{S'}(h) - \text{error}_S(h)] \right] \\ &\leq \mathbf{E}_{S,S'} \left[ \max_{h \in H} \{\text{error}_{S'}(h) - \text{error}_S(h)\} \right] \\ &= \mathbf{E}_{S,S'} \left[ \max_{h \in H} \frac{1}{m} \cdot \sum_{i=1}^m (\text{error}_{x'_i}(h) - \text{error}_{x_i}(h)) \right] \end{aligned}$$

אם  $\sigma_i = -1$  אז נבחר  $x_i \in S'$  ו- $x'_i \in S$  אם  $\sigma_i = 1$  אז נבחר  $x'_i \in S$  ו- $x_i \in S'$ . ואז:

$$\begin{aligned} \mathbf{E}[\text{MAXGAP}(S)] &\leq \mathbf{E}_S \left[ \max_{h \in H} \frac{1}{m} \cdot \sum_{i=1}^m (\text{error}_{x'_i}(h) - \text{error}_{x_i}(h)) \right] \\ &\leq \frac{2}{m} \cdot \mathbf{E}_{S,\sigma} \left[ \max_{h \in H} \sum_{i=1}^m \sigma_i \cdot \text{error}_{x_i}(h) \right] \end{aligned}$$

נשים לב כי:

$$\sigma_i \cdot (1 - 2 \cdot \text{error}_{x_i}(h)) = \sigma_i \cdot \underbrace{(c^*(x_i) \cdot h(x_i))}_{\text{error}} \simeq \sigma_i \cdot \text{error}_{x_i}(h)$$

השקילות היא בין שגיאה בטווח  $\pm 1$  לבין שגיאה בטווח  $\{0, 1\}$ .  
לכן:

$$\begin{aligned} R_D(H) &= \mathbf{E}_\sigma \left[ \max_{h \in H} \frac{1}{m} \cdot \sum_{i=1}^m \sigma_i \cdot (1 - 2 \cdot \text{error}_{x_i}(h)) \right] \\ &= 2 \cdot \mathbf{E} \left[ \max_{h \in H} \frac{1}{m} \cdot \sum_{i=1}^m \sigma_i \cdot \text{error}_{x_i}(h) \right] \end{aligned}$$

ולכן:

$$\mathbf{E}_S [\text{MAXGAP}(S)] \leq R_D(H)$$

□

## פרק 10

# שיעור 10<sup>1</sup>

סיכום השיעור אינו זמין.

---

<sup>1</sup>שיעור שהתקיים בתאריך 30.12.2012.



## **פרק 11**

# **שיעור 11<sup>1</sup>**

סיכום השיעור אינו זמין.

---

<sup>1</sup>שיעור שהתקיים בתאריך 06.01.2013.

## פרק 12

# רגרסיה<sup>1</sup>

### 12.1 הקדמה

נבחן את הדוגמה הבאה: נרצה לדעת לחזות כמה גשם ירד (בעתיד). נניח שיש מדגם  $S = \{(x_i, y_i)\}$ , כאשר  $x_i$  נדגם לפי  $D$  ו- $y_i \in \mathbb{R}$ . נסתכל על פונקציה  $f(x) = \mathbf{E}[y | x]$ . המטרה היא לבנות  $h(x)$  כך ש- $|f(x) - h(x)|$  יהיה קטן. בפועל, נרצה ש- $(f(x) - h(x))^2$  יהיה קטן. בגלל שאין לנו באמת גישה ל- $f(x)$ , כל האלגוריתמים יסתכלו על  $(y - h(x))^2$ . נרצה להוכיח שזה טוב באותה מידה.

טענה

$$\underbrace{\mathbf{E}_{x,y} [(y - h(x))^2]}_{\text{Measured error}} = \underbrace{\mathbf{E}_x [(f(x) - h(x))^2]}_{\text{Estimation error}} + \underbrace{\mathbf{E}_{x,y} [(y - f(x))^2]}_{\text{Uncertainty of the problem}}$$

הוכחה

$$\begin{aligned} \mathbf{E}_{x,y} [(y - h(x))^2] &= \mathbf{E} [(y - f(x) + f(x) - h(x))^2] \\ &= \mathbf{E} [(y - f(x))^2] + \mathbf{E} [(f(x) - h(x))^2] \\ &\quad + 2 \cdot \mathbf{E} \left[ \underbrace{(y(x) - f(x))}_{=0} \cdot (f(x) - h(x)) \right] \\ &= \mathbf{E} [(y - f(x))^2] + \mathbf{E} [(f(x) - h(x))^2] \end{aligned}$$

□

בהינתן מגדם  $S = \{(x_i, y_i)\}$  נרצה להביא למינימום את:

$$L = \frac{1}{m} \cdot \sum_{i=1}^m (y_i - h(x_i))^2$$

<sup>1</sup>סיכום לשיעור שהתקיים בתאריך 13.01.2013.

דוגמה אם  $h(x) = x \cdot w$ , אזי:

$$\begin{aligned} L &= \frac{1}{m} \cdot \sum_{i=1}^m (y_i - h(x_i))^2 \\ &= \frac{1}{m} \cdot \sum_{i=1}^m (y_i - x_i \cdot w)^2 \end{aligned}$$

נגזור:

$$\frac{d}{dw} L = \frac{1}{m} \cdot \sum_{i=1}^m (-2(y_i - x_i \cdot w) \cdot x_i) = 0$$

ולכן:

$$\frac{1}{m} \cdot \sum_{i=1}^m y_i \cdot x_i = \left( \frac{1}{m} \cdot \sum_{i=1}^m x_i^2 \right) \cdot w$$

ולכן:

$$w = \frac{\sum_{i=1}^m x_i \cdot y_i}{\sum_{i=1}^m x_i^2}$$

## 12.2 גרסיה לינארית

נניח כי  $h_w(x) = x \cdot w$ . נרצה להביא את:

$$L = \frac{1}{m} \cdot \sum_{i=1}^m (y_i - h(x_i))^2 = \frac{1}{m} \cdot \sum_{i=1}^m (y_i - x_i \cdot w)^2$$

למינימום.

נסתכל על זה כמטריצות. נסמן:

$$W = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix}$$

$$X = \begin{pmatrix} -x_1- \\ -x_2- \\ \vdots \\ -x_m- \end{pmatrix}$$

אז:

$$X \cdot W - y = \begin{pmatrix} x_1 \cdot W - y_1 \\ x_2 \cdot W - y_2 \\ \vdots \\ x_m \cdot W - y_m \end{pmatrix}$$

כאשר  $W \in \mathbb{R}^n$ ,  $x_i \in \mathbb{R}^n$  ו- $y_i \in \mathbb{R}$ .  
אזי:

$$\begin{aligned} \min_w L &= \min_w (X \cdot W - y)^T \cdot (X \cdot W - y) \\ &= \min_w (X \cdot W)^T \cdot (X \cdot W) - 2 \cdot (X \cdot W) \cdot y + y \cdot y \end{aligned}$$

הגרדיאנט:

$$\begin{aligned} \nabla_W L &= 2 \cdot X^T \cdot X \cdot W - 2 \cdot X^T \cdot y = 0 \\ w &= (X^T \cdot X)^{-1} \cdot X^T \cdot y \end{aligned}$$

ראינו כי לרגסיה לינארית יש פתרון סגור. ניתן להכליל את הפתרון גם לרגסיה מסדר יותר גבוה, בעזרת שימוש ב-Kernel של מרחב וקטורי כלשהו.

### 12.2.1 חוסר יציבות הפתרון

אחת הבעיות המשמעותיות של הרגסיה הלינארית היא חוסר יציבות הפתרון. היינו מצפים ששני קלטים מאוד דומים יתנו תוצאות דומות. למשל,  $x_1 = (1, 0)$ ,  $x_2 = (1, \varepsilon)$ ,  $y_1 = 1 = y_2$ . אזי:

$$\begin{aligned} L &= \min_w \left[ \frac{1}{2} \cdot (w_1 - 1)^2 + \frac{1}{2} \cdot (w_1 + \varepsilon \cdot w_2 - 1)^2 \right] \\ \frac{d}{dw_1} L &= (w_1 - 1) + (w_1 + \varepsilon \cdot w_2 - 1) = 0 \\ \frac{d}{dw_2} L &= (w_1 + \varepsilon \cdot w_2 - 1) \cdot \varepsilon = 0 \\ w &= (w_1, w_2) = (1, 0) \end{aligned}$$

נשנה את  $y_1$  להיות  $1 + \varepsilon$ . עכשיו:

$$\begin{aligned} L &= \min_w \left[ \frac{1}{2} \cdot (w_1 - 1 - \varepsilon)^2 + \frac{1}{2} \cdot (w_1 + \varepsilon \cdot w_2 - 1)^2 \right] \\ \frac{d}{dw_1} L &= (w_1 - 1 - \varepsilon) + (w_1 + \varepsilon \cdot w_2 - 1) = 0 \\ \frac{d}{dw_2} L &= \varepsilon \cdot (w_1 + \varepsilon \cdot w_2 - 1) = 0 \end{aligned}$$

הפתרון יהיה  $w_1 = 1 + \varepsilon$ ,  $w_2 = -1$ .

## 12.3 רגולריזציה

### 12.3.1 Ridge Regression

נגדיר *Ridge Regression*: נרצה להביא למינימום את:

$$\frac{1}{2} \cdot \sum_{i=1}^m \underbrace{(w \cdot x_i - y_i)^2}_{sQ=Quadratic\ Loss} + \frac{\lambda}{2} \cdot \|w\|^2$$

כאשר  $\frac{\lambda}{2} \cdot \|w\|^2$  הוא המשקל של  $w$ , ו- $\lambda$  הוא קבוע ל-trade off בין שני החלקים ברגרסיה. בכתוב מטריציוני:

$$\begin{aligned} L &= \frac{1}{2} \cdot (X \cdot w - y)^T \cdot (X \cdot w - y) + \frac{\lambda}{2} \cdot W^T \cdot W \\ \nabla_W L &= X^T \cdot X \cdot W - X^T \cdot y + \lambda \cdot W = 0 \\ \underbrace{(X^T \cdot X + \lambda \cdot \mathbb{I})}_{\text{Eign values} \geq \lambda} \cdot W &= X^T \cdot y \\ W &= \underbrace{(X^T \cdot X + \lambda \cdot \mathbb{I})^{-1}}_{\text{Eign values} \leq 1/\lambda} \cdot X^T \cdot y \end{aligned}$$

### 12.3.2 Lasso Regression

יש עוד דרך לייצב את הרגרסיה: *Lasso Regression*. נרצה להביא למינימום את  $SQ - \|w\|_1 \cdot \lambda$ . קשה לנמק את זה מבחינה אנליטית, ולכן אין פתרון סגור. לעומת זאת, קל לפתור את זה בתוכנה. הפתרון יהיה דליל במספר המשקולות. אי אפשר להוכיח את זה בצורה מתמטית, אבל אפשר להבין את זה על סמך ניסיון.

### 12.3.3 חסם הכללה ל-Ridge Regression (או משהו שדומה לו)

נחפש את:

$$\min_w \sum_{i=1}^m \psi_i^2$$

כאשר:

$$\psi_i^2 = w \cdot x_i - y_i$$

ונדרוש:

$$\|w\|_2^2 \leq \Lambda^2$$

נניח גם כי:

$$\|x_i\| \leq R$$

**למה (Ledoux-Talag)** נניח כי  $\Phi_i: \mathbb{R} \rightarrow \mathbb{R}$  היא  $M$ -ליפשיץ, כלומר:

$$\forall a, b. |\Phi_i(a) - \Phi_i(b)| \leq M \cdot |a - b|$$

אז:

$$\begin{aligned} \hat{R}_S(\Phi \circ H) &= \mathbf{E}_\sigma \left[ \sup_{h \in H} \frac{1}{m} \cdot \sum_{i=1}^m \sigma_i \cdot \Phi(h(x_i)) \right] \\ &\leq M \cdot \mathbf{E}_\sigma \left[ \sum_{h \in H} \left| \frac{1}{m} \cdot \sum_{i=1}^m \sigma_i \cdot h(x_i) \right| \right] \\ &= M \cdot \hat{R}_S(H) \end{aligned}$$

במקרה שלנו, נניח כי  $\psi_i = w \cdot x_i - y_i$  חסום על ידי  $M$ . אזי  $\Phi(\psi) = \psi^2$  היא  $2M$ -ליפשיץ. מהלמה, נקבל כי:

$$\hat{R}_S(SQ) \leq 2 \cdot M \cdot \hat{R}_S(\text{linear}) \leq 2 \cdot M \cdot \frac{R \cdot \Lambda}{\sqrt{m}}$$

מכאן, בהסתברות  $1 - \delta$ :

$$\mathbf{E}[SQ] \leq SQ(S) + 2 \cdot \hat{R}_S(SQ) + O\left(\sqrt{\frac{\log 1/\delta}{m}}\right)$$

#### 12.3.4 נקודת מבט בייסיאנית

נגדיר:  $y = w \cdot x + \text{Noise}$ , כאשר  $\text{Noise} \sim N(0, \sigma)$  ו- $w \sim N(0, \sigma_0)$ . לפי חוק Bayes:

$$\Pr[w | \{(x_i, y_i)\}] = \frac{\Pr[\{(x_i, y_i)\} | w] \cdot \Pr[w]}{\Pr[\{(x_i, y_i)\}]}$$

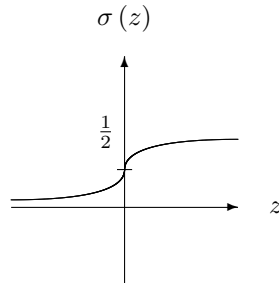
לפי שיטת Maximum Likelihood:

$$\begin{aligned} w_{\text{ML}} &= \arg \max_w \Pr[\{(x_i, y_i)\} | w] \\ &= \arg \max_w \prod_{i=1}^m \Pr[(x_i, y_i) | w] \\ &= \arg \max_w e^{-\sum_{i=1}^m \frac{(y_i - w \cdot x_i)^2}{2 \cdot \sigma^2}} \\ &= \arg \min_w \sum_{i=1}^m \frac{1}{2 \cdot \sigma^2} \cdot (y_i - w \cdot x_i)^2 \end{aligned}$$

קיבלנו מינימום של שגיאה ריבועית, כלומר גרסיה לינארית. לפי שיטת MAP:

$$\begin{aligned} w_{\text{MAP}} &= \arg \max_w \Pr[\{(x_i, y_i)\} | w] \cdot \Pr[w] \\ &= \arg \max_w e^{-\sum_{i=1}^m \frac{(y_i - w \cdot x_i)^2}{2 \cdot \sigma^2}} \cdot e^{-\frac{\|w\|_2^2}{2 \cdot \sigma_0^2}} \\ &= \arg \min_w \underbrace{\sum_{i=1}^m \frac{(y_i - x_i \cdot w)^2}{2 \cdot \sigma^2}}_{\text{Quadratic error}} + \underbrace{\frac{\|w\|_2^2}{2 \cdot \sigma_0^2}}_{\text{Weight of } w} \end{aligned}$$

קיבלנו Ridge Regression עם  $\lambda = \frac{\sigma^2}{\sigma_0^2}$ . התרגיל הזה הראה לנו שאפשר לגזור את הרגרסיות גם מתוך עולם בייסיאני.

איור 12.1: תיאור של  $\sigma(z)$ 

## Logistic Regression 12.4

נחזור לבעיות סיווג: נניח כי  $Y = \{0, 1\}$  או  $Y = \{+1, -1\}$ . עדיין יש משמעות ללדבר על  $\Pr[Y = 1 | X]$ . אזי:

$$\mathbf{E} \left[ (w \cdot x_i - y_i)^2 \right] = \mathbf{E} [w \cdot x - \Pr[y = 1 | x]]$$

הרגרסיה תבנה מודל שבו

$$\Pr[Y = 1 | X] \simeq w \cdot x$$

עד כאן הכל עובד כמו שצריך. אבל לא מובטח ש- $w \cdot x \in \{0, 1\}$ . נרצה להחזיר את הפתרון שלנו חזרה לעולם של בעיות הסיווג. לכן, נאמר כי

$$\Pr[Y = 1 | X] = \sigma(w \cdot x)$$

כאשר

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

נשים לב לתכונה הבאה:  $\sigma(z) + \sigma(-z) = 1$ . למה זה מתקיים?

$$\begin{aligned} \sigma(z) + \sigma(-z) &= \frac{1}{1 + e^z} + \frac{1}{1 + e^{-z}} \\ &= \frac{\frac{1+e^z}{1+e^{-z}} + 1}{1 + e^z} \\ &= \frac{1 + e^z + 1 + e^{-z}}{(1 + e^z) \cdot (1 + e^{-z})} \\ &= \frac{2 + e^z + e^{-z}}{1 + e^z + e^{-z} + 1} \\ &= \frac{2 + e^z + e^{-z}}{2 + e^z + e^{-z}} = 1 \end{aligned}$$

למה זה טוב? אם  $\Pr[Y = 1 | X] = \sigma(w \cdot x)$  אזי:

$$\begin{aligned} 1 - \Pr[Y = 1 | X] &= \Pr[Y = -1 | X] \\ &= \sigma(-w \cdot x) = 1 - \sigma(w \cdot x) \end{aligned}$$

נרצה למצוא  $w$  שיתן לנו את

$$\min_w \sum_{i=1}^m (\sigma(y_i \cdot w \cdot x_i) - y_i)^2$$

נחפש מינימום מקומי. אזי:

$$\begin{aligned} w &= \arg \max_w \prod_{i=1}^m \Pr[y_i | x_i \cdot w] \\ &= \arg \max_w \sum_{i=1}^m \log \Pr[y_i | x_i \cdot w] \\ &= \arg \max_w \sum_{i=1}^m \log \sigma(y_i \cdot x_i \cdot w) \\ &= \arg \min_w \sum_{i=1}^m -\log \sigma(y_i \cdot x_i \cdot w) \\ &= \arg \min_w \sum_{i=1}^m \log(1 + e^{y_i \cdot x_i \cdot w}) \end{aligned}$$

איך נשערך את ההסתברויות הבאות:

$$\begin{aligned} \Pr[Y = 1 | X] &= \sigma(w \cdot x_i) \\ \Pr[Y = -1 | X] &= \sigma(-w \cdot x_i) \end{aligned}$$

נהפוך את זה למסווג:  $\sigma(w \cdot x) \geq 1/2$  או  $w \cdot x \geq 0$ . קיבלנו בחזרה מפריד לינארי.



## פרק 13

# 1 Model Selection

### 13.1 הקדמה

בהינתן מספר דוגמאות נתון, נרצה למצוא השערה טובה ביותר. נניח כי  $c^*$  היא פונקציית המטרה ו- $H$  היא מחלקת ההשערות. לשם פשטות (לפחות לאינטואיציה בהתחלה), נניח כי  $c^* \in H$ . נניח כי  $H$  בת-מניה:  $H = \{h_1, \dots, h_i, \dots\}$ . אלגוריתם 1.13 הוא אלגוריתם Model Selection. ננתח אותו: ראשית, האלגוריתם תמיד עוצר. מכיוון ש- $c^* \in H$ , אזי  $\exists j. h_j = c^*$ . אם נגיע לפאזה  $j$ , אז בוודאי נעצור, כי  $c^* = h_j$  תהיה עקבית. לכן, מספר הפאזות חסום על ידי  $j$ . כמו כן, אם האלגוריתם עוצר ומחזיר את  $h_i$ , אזי מתקיים:

$$\Pr [h_i \text{ is consistent} \mid \text{error}(h_i) > \varepsilon] \leq \delta_i$$

למה? ההסתברות שנעצור היא ההסתברות שב- $m_i$  לא ראינו שגיאות. ההסתברות לכך היא:

$$(1 - \text{error}(h_i))^{m_i} \leq (1 - \varepsilon)^{m_i} \leq e^{-\varepsilon \cdot m_i} = \delta_i$$

נבחן את פרמטר הביטחון:

$$\delta = \sum_{i=1}^{\infty} \delta_i$$

למשל,  $\delta_i = \frac{1}{2^i}$  או  $\delta_i = \frac{1}{i^2}$  נותנים  $\delta = 1$  או  $\delta = \frac{\pi^2}{6}$  בהתאמה. לרוב נגדיר  $\delta_i = \frac{\delta}{i^2} \cdot \frac{6}{\pi^2}$  ונקבל פרמטר ביטחון  $\delta$ .

<sup>1</sup>סיכום לשיעור שהתקיים ב-20.01.2013.

---

#### אלגוריתם 1.13 Model Selection

נרוץ בפאזות. בפאזה ה- $i$ :

• נדגום  $m_i = \frac{1}{\varepsilon} \cdot \ln \frac{1}{\delta_i}$  דוגמאות.

– אם  $h_i$  עקבי, נחזיר את  $h_i$ .

– אחרת, נמשיך לפאזה הבאה.

---

שגיאות	החלפות	
		דגימה
		0 7
		1 3
		7 0

טבלה 13.1: דוגמאות לדגימות והשערות עליהן

נבחן את מספר הדוגמאות: בלי מחזור של דוגמאות, נשתמש ב- $\sum_{i=1}^j m_i$  דוגמאות. עם מחזור:

$$\max_{1 \leq i \leq j} \{m_i\} = m_j$$

כאשר:

$$m_j = \frac{1}{\varepsilon} \cdot \ln \left( \frac{j^2}{\delta} \cdot \frac{6}{\pi^2} \right)$$

כלומר, מספר הדוגמאות תלוי באינדקס של פונקציית המטרה.

### 13.1.1 דוגמה

ניתן לראות בטבלה 13.1 דוגמה לדגימות והשערות עליהן. היינו מעדיפים לקבל את האפשרות השנייה: האפשרות הראשונה מפספסת יותר מדי, ואילו האפשרות השלישית יותר מדי מותאמת לדגימה. האפשרות השנייה היא משהו טוב ביניהן.

### 13.1.2 המודל

נתונה משפחת השערות  $H$ , כך ש- $H_1 \subseteq H_2 \subseteq H_3 \subseteq \dots \subseteq H_i \subseteq \dots \subseteq H$ . כמו כן:

$$H = \bigcup_{i=1}^{\infty} H_i$$

הסיבוכיות של  $H_i$ :  $\text{VCdim}(H_i) = i$ ,  $|H_i| = 2^{i-1}$ ,  $|H_1| = 1$ . נניח כי פונקציית המטרה  $c^* \notin H$  מקיימת  $c^* \notin H$ . נסמן:

$$\varepsilon(h) = \Pr[h \neq c^*]$$

$$\varepsilon_i = \min_{h \in H_i} \varepsilon(h)$$

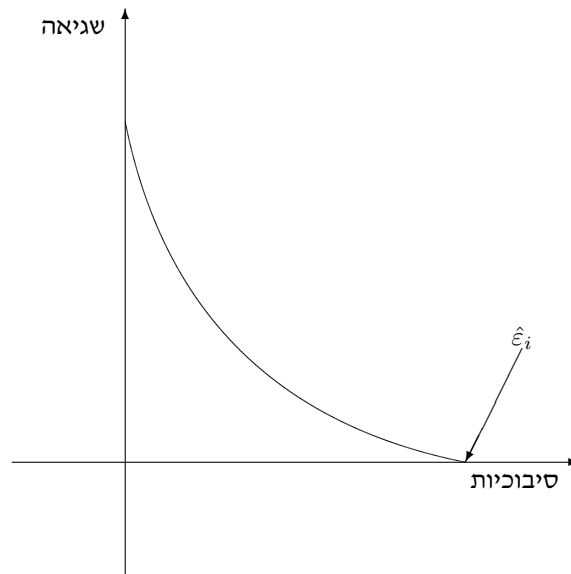
בגלל ש- $H_i \subseteq H_{i+1}$ , אזי  $\varepsilon_i \geq \varepsilon_{i+1}$ , ולכן נגדיר:

$$\varepsilon^* = \inf_i \varepsilon_i$$

כמו כן, נגדיר:

$$\hat{\varepsilon}(h) = \frac{1}{m} \cdot \sum_{x \in S} \mathbf{I}(h(x) \neq c^*(x))$$

$$\hat{\varepsilon}_i = \min_{h \in H_i} \hat{\varepsilon}(h)$$



איור 13.1: השגיאה כפונקציה של ההסתברות כאשר מסתכלים על  $\hat{\epsilon}_i$ .

מכיוון ש- $\text{VCdim}(H) = \infty$ , אזי כל  $m$  נקודות יכולות לקבל את כל הסיווגים. לכן,  $\hat{\epsilon}_m = 0$  היינו יכולים לבחור את  $g = \arg \min_h \hat{\epsilon}(h)$ , אבל תמיד קיים  $h \in H_m$  כך ש- $\hat{\epsilon}(h) = 0$ , ולכן אין בהכרח משמעות ל- $g$  (ניתן לראות את זה באיור 13.1). אפשר לבחור:  $g^* = \arg \min_h \{\hat{\epsilon}(h) + p(h)\}$ . נגדיר:

$$d(h) = \min \{i \mid h \in H_i\}$$

### Structural Risk Minimization 13.2

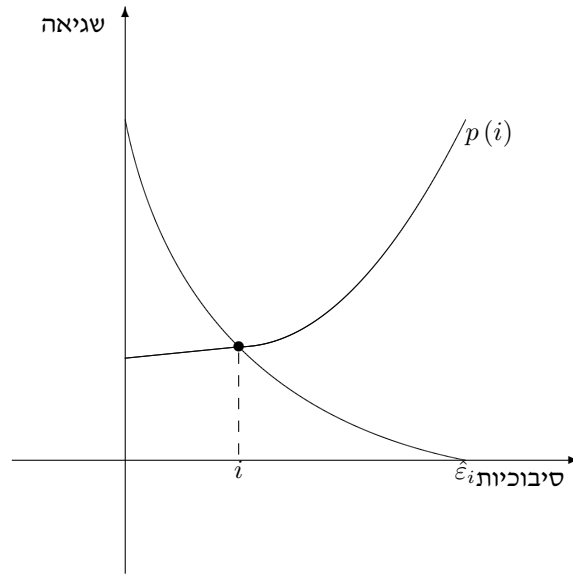
ב-Structural Risk Minimization (או בקיצור SRM), נבחר:

$$\text{Penalty}(h) = p(h) = \sqrt{\frac{2 \cdot d(h) \cdot \ln 2 + \ln 1/\delta}{m}}$$

כאשר  $m$  הוא מספר הדוגמאות ו- $\delta$  הוא פרמטר הבטחון. נבחר:

$$g^* = \arg \min_{h \in H} \hat{\epsilon}(h) + \sqrt{\frac{2 \cdot d(h) \cdot \ln 2 + \ln 1/\delta}{m}}$$

נסמן:  $h^* = \arg \min_{h \in H} \epsilon(h)$



איור 13.2: השגיאה כפונקציה של ההסתברות כאשר מוסיפים את Penalty להשערות ארוכות

**משפט** בהסתברות  $1 - \delta$ :

$$\varepsilon(g^*) \leq \varepsilon(h^*) + 2 \cdot \text{Penalty}(h)$$

**הוכחה** ניזכר בטענה כי  $\Pr[|\varepsilon(h) - \hat{\varepsilon}(h)| \geq \lambda] \leq 2 \cdot e^{-\lambda^2 \cdot m}$  (נובעת ישירות מ-Chernoff). נבחר:

$$g_i = \arg \min_{h \in H_i} \hat{\varepsilon}(h)$$

אז:

$$\begin{aligned} \Pr[|\varepsilon(g_i) - \hat{\varepsilon}(h_i)| \geq \lambda_i] &\leq \Pr[\exists h \in H_i. |\varepsilon(h) - \hat{\varepsilon}(h)| \geq \lambda_i] \\ &\leq |H_i| \cdot 2 \cdot e^{-\lambda_i^2 \cdot m} = 2^i \cdot e^{-\lambda_i^2 \cdot m} = \delta_i \end{aligned}$$

נקבע:

$$\lambda_i = \sqrt{\frac{i \cdot \ln 2 + \ln 2^i / \delta}{m}}$$

נחשב את ההסתברות עבור כל ה- $i$ ים:

$$\begin{aligned} \Pr[\forall i. \forall h \in H_i. |\varepsilon(h) - \hat{\varepsilon}(h)| \leq \lambda_i] &= 1 - \Pr[\exists i. \exists h \in H_i. |\varepsilon(h) - \hat{\varepsilon}(h)| \geq \lambda_i] \\ &= 1 - \sum_{i=1}^{\infty} \delta_i = 1 - \delta \end{aligned}$$

אפשר לבחור  $\delta_i = \frac{\delta}{2^i}$  או  $\delta_i = \frac{\delta}{i^2} \cdot \frac{6}{\pi^2}$ . לכן, בהסתברות  $1 - \delta$ , לכל  $h \in H$  מתקיים  
 $|\varepsilon(h) - \hat{\varepsilon}(h)| \leq \lambda_{d(h)}$   
 נציב את  $\delta_i$  ב- $\lambda_i$ :

$$\lambda_i = \sqrt{\frac{i \cdot \ln 2 + \ln^2 i / \delta}{m}} = \sqrt{\frac{2 \cdot i \cdot \ln 2 + \ln 1 / \delta}{m}}$$

נניח כעת כי  $d(h^*) = j$  ו- $d(g^*) = i$  אזי:

$$\begin{aligned} \hat{\varepsilon}(h^*) &\leq \varepsilon(h^*) + \lambda_i \\ \varepsilon(g^*) - \lambda_j &\leq \hat{\varepsilon}(g^*) \end{aligned}$$

כמו כן, מבחירת  $g^*$ :

$$\begin{aligned} \hat{\varepsilon}(g^*) + p_j &\leq \hat{\varepsilon}(h^*) + p_i \\ \downarrow \\ \varepsilon(g^*) - \lambda_j + p_j &\leq \varepsilon(h^*) + \lambda_j + p_j \end{aligned}$$

נשים לב כי  $\lambda_k = p_k$ . לכן:  $\varepsilon(g^*) \leq \varepsilon(h^*) + 2 \cdot p_i$ .  $\square$

### (CV) Cross Validation 13.3

נשים לב: ב-SRM בשלב א', מכל  $H_i$  בוחרים  $g_i$ , ובשלב ב' בוחרים בין  $g_i$  בעזרת  $p_i$ . ב-CV, שלב א' זהה, ובשלב ב' מערבים דוגמאות חדשות. נתון מדגם  $S$  בגודל  $m$ . נחלק ל- $S_1$  ו- $S_2$ , כאשר  $|S_1| = (1 - \gamma) \cdot m$  ו- $|S_2| = \gamma \cdot m$ .  $S_1$  ישמש את שלב א', ו- $S_2$  ישמש את שלב ב'.

**שלב א'** לכל  $H_i$  נבחר  $g_i = \arg \min_{h \in H_i} \hat{\varepsilon}_1(h)$  כאשר  $\hat{\varepsilon}_1$  פועל על  $S_1$ . נגדיר:

$$G = \{g_i \mid \hat{\varepsilon}_{i,1}(g_i) < \hat{\varepsilon}_{i-1,1}(g_i)\}$$

נשים לב כי  $|G| \leq m$ .

**שלב ב'** נבחר מתוך  $G$  בעזרת  $S_2$ . כלומר:

$$g^* = \arg \min_{g \in G} \hat{\varepsilon}_2(g)$$

כאשר  $\hat{\varepsilon}_2$  פועל על  $S_2$ . יהי  $\mathcal{A}$  אלגוריתם כללי שבוחר מתוך  $G$ .

#### משפט (Cross-Validation)

$$\varepsilon_{CV}(m) \leq \varepsilon_{\mathcal{A}}((1 - \gamma) \cdot m) + 2 \cdot \sqrt{\frac{\ln \frac{2 \cdot m}{\delta}}{m}}$$

**הוכחה** נניח כי  $g_i$  נבחר על ידי  $\mathcal{CV}$ , ונשווה אותו ל- $g_k$  כללי. נרצה להראות:

$$\varepsilon(g_i) \leq \varepsilon(g_k) + 2 \cdot \sqrt{\frac{\ln \frac{2 \cdot m}{\delta}}{m}}$$

טענת ההכללה:

$$\begin{aligned} \forall g \in G. \Pr[|\varepsilon(g) - \hat{\varepsilon}(g)| \geq \lambda] &\leq 2 \cdot e^{-\lambda^2 \cdot \gamma \cdot m} \\ &\Downarrow \\ \Pr[\exists g \in G. |\varepsilon(g) - \hat{\varepsilon}(g)| \geq \lambda] &\leq 2 \cdot |G| \cdot e^{-\lambda^2 \cdot \gamma \cdot m} \\ &\leq 2 \cdot m \cdot e^{-\lambda^2 \cdot \gamma \cdot m} = \delta \\ &\Downarrow \\ \lambda &= \sqrt{\frac{\ln \frac{2 \cdot m}{\delta}}{\gamma \cdot m}} \\ &\Downarrow \\ \Pr[\exists g \in G. |\varepsilon(g) - \hat{\varepsilon}(g)| \geq \lambda] &\leq \delta \end{aligned}$$

לכן, בהסתברות  $1 - \delta$ :

$$\begin{aligned} \varepsilon(g_i) - \lambda &\leq \hat{\varepsilon}_2(g_i) \\ \hat{\varepsilon}_2(g_k) &\leq \varepsilon(g_k) + \lambda \\ \hat{\varepsilon}_2(g_i) &\leq \hat{\varepsilon}_2(g_k) \end{aligned}$$

לכן,  $\varepsilon(g_i) - \lambda \leq \varepsilon(g_k) + \lambda$ , ולכן  $\lambda \leq \varepsilon(g_k) - \varepsilon(g_i) + 2 \cdot \lambda$ .  
מההגדרה,  $\mathbf{E}[\varepsilon(g_i)] = \varepsilon_{\mathcal{CV}}(m)$ . כמו כן:

$$\mathbf{E} \left[ \min_{g_k \in G} \varepsilon(g_k) \right] \leq \varepsilon_{\mathcal{A}}((1 - \gamma) \cdot m)$$

לכן:

$$\varepsilon_{\mathcal{CV}}(m) \leq \varepsilon_{\mathcal{A}}((1 - \gamma) \cdot m) + 2 \cdot \sqrt{\frac{\ln \frac{2 \cdot m}{\delta}}{m}}$$

□

10-fold Cross Validation מחלק ל- $\{S_i\}_{i=1}^{10}$  כאשר  $|S_i| = \frac{m}{10}$ . האלגוריתם משתמש ב- $S_i$  ימים, ובודק על הנתר.

### 13.4 (MDL) Minimum Description Length

בשלב הראשון נמצא קידוד ל- $h$ , ובשלב השני נתקן. אם יש הרבה תיקונים, עדיף למצוא קידוד אחר ל- $h$ .  
לכן, נחפש:

$$g^* = \arg \min_h \left\{ \underbrace{\text{size}(h)}_{p(h)} + \underbrace{\text{size}(\text{corrections})}_{f(\varepsilon(h))} \right\}$$

**13.4.1 בעזרת MAP**

לפי ההגדרה:

$$\Pr[h | D] = \frac{\Pr[D | h] \cdot \Pr[h]}{\Pr[D]}$$

לכן:

$$\begin{aligned} g^* &= \arg \max_h \Pr[D | h] \cdot \Pr[h] \\ &= \arg \max_h \log \Pr[D | h] + \log \Pr[h] \\ &= \arg \min_h \frac{1}{\log \Pr[D | h]} + \frac{1}{\log \Pr[h]} \end{aligned}$$