מטריצת אימון: $X \in \mathbb{R}^{d imes m}$ ששורותיה הן תכונות, ועמודותיה הן עבור $y_i = f(x_i) + z_i$ בעל ערך $y \in \mathbb{R}^m$ עבור $y_i = y_i$. רעש כלשהו z_i

 $:\!\ell \colon \mathbb{R}^d o \mathbb{R}$ עבור פונקציית מחיר: Empirical Risk

$$ER(h) = \frac{1}{m} \sum_{i=1}^{m} \ell(h(x_i), y_i)$$

i=1		
Squared Loss (RSS)	$\sum_{i=1}^{m} (h(x_i) - y_i)^2 = \ h(X) - y\ _2^2$	
Absolute Value Loss	$\sum_{i=1}^{m} h(x_i) - y_i = h(X) - y _1$	
0-1 Loss	$\frac{1}{m}\sum_{i=1}^{m} 1_{[h(x)\neq y]}$	

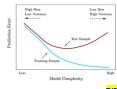
Maximum Likelihood

 $\mathcal{D}(\Theta)$ עם פרמטר: בהינתו התפלגות בהינתו פרמטר: יודגימות או מגדירים את הנראות: $x_1,\dots,x_m \stackrel{iid}{\sim} \mathcal{D}(\Theta)$ ודגימות $L(\Theta) = L(\Theta|x_1,\dots,x_m) = \Pr_{\Gamma}[x_1,\dots,x_m|\Theta]$

X =יהי (Maximum Likelihood Estimator ר"ת) MLE המרבית: קבוצת אז הנראות קבוצת דגימות, אז הנראות המרבית: $x_1, \dots, x_m \overset{iid}{\sim} \mathcal{D}(\Theta)$ $\widehat{\Theta}_{MLE} = \arg \max_{\Omega} L(\Theta|X)$

Bias: הסטייה של המודל מנקודות המדגם (שגיאת ההכללה). . ככל שקבוצת ההיפותזות גדולה יותר, כך הסטייה קטנה. נובעת מ-

Variance: השונות של המודל מהאמת (מרדף אחרי הרעש). ככל שקבוצת ההיפותזות גדולה יותר, כך השונות גדלה. נובעת מ-



 $\mathcal{H}_{reg} = \{h(x) | h(x) = \langle w, \binom{1}{x} \rangle, w \in \mathbb{R}^{d+1} \}$ (נקראים המשקולות נקרא $w_1, ..., w_d$.intercept נקרא נקרא w_0)

 $X(y - X^T w) = 0$ שקול ל- RSS-מזעור ה-**Normal Equations טענה:** התאים הבאים שקולים: 1. קיים פתרון יחיד למערכת. $.\sigma_{d+1}(X)>0$.4. הפיכה XX^T .3 .dim $\left(Ker(X^T)\right)=0$.2 $\widehat{w} = (XX^T)^{-1}Xy$ (בעל נורמה מינימלית). $\widehat{w} = (XX^T)^{-1}Xy$ טענות מההוכחה: 1. $\operatorname{Ker}(X^{\mathrm{T}}) = \operatorname{Ker}(X^{\mathrm{T}})$. עבור מטריצה

מערכת לא $y=X^Tw$ תהי $\mathbf{3}.Im(A^T)=Ker(A)^\perp$ מערכת לא ריבועית, ביועית, לא הפיך. למערכת אינסוף פתרונות אם"ם אומוגנית, ונניח X^T לא הפיך. למערכת אינסוף יש אינסוף פתרונות או $XX^Tw=Xy$ למערכת **4** y $\perp Ker(X)$ XX^T פתרון יחיד (אם XX^T הפיכה).

משפט: $y^{\dagger}(X^T)$ הוא תמיד פתרון למשוואה הנורמלית. הגדרה: אם X^T כולל עמודה (פיטצ'ר) הנפרש ע"י יתר העמודות נאמר שהוא כמעט סינגולרי (=לא הפיך). במקרה זה ניקח פתרון:

$$X^{\dagger,\epsilon} = V\Sigma^{\dagger,\epsilon}U^T \text{ s.t } [\Sigma^{\dagger,\epsilon}]_{ii} = \begin{cases} \frac{1}{\epsilon_i} & \sigma_i > \epsilon \\ 0 & otherwise \end{cases}$$

אז $(y_i \overset{iid}{\sim} N(x_i^T w, \sigma^2)$ כלומר $z_i \overset{iid}{\sim} N(0, \sigma^2)$ אז :MLE

$$p(y|w) = \prod_{i=1}^{m} \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{\left(x_i^{T_w} - y_i\right)^2}{2\sigma^2}}$$

$$\frac{1}{2\sigma^2} \prod_{i=1}^{m} \frac{-\frac{\left(x_i^{T_w} - y_i\right)^2}{2\sigma^2}}{2\sigma^2} I(x_i|x_i)$$

 $= \frac{1}{(\sqrt{2\pi\sigma^2})^m} \prod_{i=1}^m e^{-\frac{(x_i^T w - y_i)^2}{2\sigma^2}} = L(w|y)$

 $\widehat{w} = \arg \max_{w} L(w|y) \stackrel{by log}{=} \arg \min_{w} \sum_{i=1}^{m} (x_i^T w - y_i)^2$

-מתקיים: Bias-Variance Tradeoff

 $MSE = \mathbb{E}[(\hat{y} - y)^2] = Var[\hat{y}] + bias^2[\hat{y}]$

סיווג שגיאות: שגיאת <u>Type-l</u> (או FP) זו השגיאה החמורה.

שגיאת **Type-II** (או FN).

True	TP, Recall, TPR:	TN, Specificity:	Accuracy: TP + TN
	Positive	Negative	Total
False	FP, FPR: FP	FN, FNR: FN	Error Rate: FP + FN
	Negative	Positive	Total
	$\frac{TP}{TP + FP}$	Negative predictive value: TN $TN + FN$	$F_1 \\ = 2 \cdot \frac{Recall \cdot Precision}{Recall + Precision}$
			11-16

$\mathcal{H}_{half} = \{x \mapsto sign(\langle x, w \rangle + b)\}$

תחת ההנחה שקיים מפריד לינארי, נגדיר את הסיכון להיות

.ERM-ונמזער לפי עקרון וונ $L_S(h_w) = \sum_{i=1}^m 1_{\{y_i\cdot (x_i,w)<0\}}$ arg min 0

 $s.t \quad diag(y) \cdot X^T w \geq 1$

Perceptron Algorithm:

Initialize: $w = \vec{0}$ While $\exists i \text{ s.t } y_i \cdot \langle x_i, w \rangle \leq 0$:

 $w \leftarrow w + y_i \cdot x_i$

 $\mathcal{H}_{SVM} = \{x \mapsto sign(\langle x, w \rangle + b)\}$ $D\left(\binom{w}{b}, x_i\right) = \min_{v: \langle v, w \rangle + b = 0} \|x_i - v\|$ מרחק: $M\left(\binom{w}{b}, S\right) = \min_{w \in S} D\left(\binom{w}{b}, x_i\right)$: שוליים

 $Hard - SVM = \arg\min_{w,h} ||w||^2$ $\begin{array}{ccc} & & & & \\ s.t & y_i \cdot (\langle x_i, w \rangle + b) \geq 1 \end{array} \quad (margin = \frac{2}{\|w\|})$ $Soft - SVM = \arg\min_{w \mid h} \lambda \cdot ||w||^2 + \frac{1}{m} \sum_{i=1}^{m} \xi_i$

s.t $y_i \cdot (\langle x_i, w \rangle + b) \ge 1 - \xi_i$ $\xi_i \ge 0$

 $\Leftrightarrow \arg\min_{w} \lambda \cdot \|w\|^2 + \frac{1}{m} \sum_{i=1}^{m} \underbrace{\max\{1 - y_i \cdot (\langle x_i, w \rangle + b), 0\}}_{}$

מרחב מכפלה F מרחב (באשר $\psi \colon \mathcal{X} \to F$ מציאת העתקה (באשר וציאת ביאת) פנימית שלם). $w^* = \arg\min \lambda \cdot ||w||^2 + f(\langle w, \psi(x_1) \rangle, ..., \langle w, \psi(x_m) \rangle)$

טענה: קיים $\alpha \in \mathbb{R}^m$ עבורו $\alpha \in \mathbb{R}^m$. במקרה במקרה מענה: ואז: $K(x_i, x) = \langle \psi(x_i), \psi(x) \rangle$

 $\langle w^*, \psi(x) \rangle = \sum \alpha_i K(x_i, x)$

Bayes Optimal Classifier

 $L_{\mathcal{D}}(h) =$ בהינתן התפלגות \mathcal{D} מעל $\mathcal{X} imes \mathcal{Y}$, עבור הפסד :נגדיר , $\mathbb{E}_{(x,y)\sim\mathcal{D}}[h(x)\neq y]$

 $f_{\mathcal{D}} = \arg \max_{y \in \mathcal{T}} \Pr[Y = y | X = x]$

 $L_{\mathcal{D}}(f_{\mathcal{D}}) \leq L_{\mathcal{D}}(h)$ מתקיים $h \colon \mathcal{X} o \mathcal{Y}$ טענה: עבור היפותזה נהור: $x|y \sim \mathcal{N}(\mu_y, \Sigma)$ והתפלגות $\mathcal{X} = \mathbb{R}^d$ נכלומר: LDA והתפלגות 1 $f(x|y) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} e^{-\frac{1}{2}(x-\mu_y)^T \Sigma^{-1}(x-\mu_y)}$

 $h_{\mathcal{D}}(x) = \arg \max_{y} x^{T} \Sigma^{-1} \mu_{y} - \frac{1}{2} \mu_{y}^{T} \Sigma^{-1} \mu_{y} + \ln(\Pr[y])$

Logistic Regression

 $\phi \colon \mathbb{R} \to \underbrace{\sum_{i=1}^{n} N(0, \sigma^2)}_{i \in N(0, \sigma^2)}$, פונקציה y = Xw + z[0,1] (מונוטונית עולה, חלקה, חח"ע מהישר הממשי לקטע $w \in \gamma_i = \phi(\langle x_i, w \rangle)$ עבור $y_i \sim Ber(p_i)$ נניח ((0,1).

 $\pi(x)=rac{e^x}{1+e^x}=rac{1}{1+e^{-x}}$ נשתמש בפונקציה הלוגיסטית $\mathcal{H}_{logistic}^d=\{x\mapsto \pi(\langle x,w
angle)\}$

 $L(w|y) = \prod p_i(w)^{y_i} \left(1-p_i(w)\right)^{1-y_i}$ התועלת, ובעזרת $\widehat{w} = \arg\max_{w \in \mathbb{R}^{d+1}} \sum \left[y_i \langle x_i, w \rangle - \log \left(1 + e^{\langle x_i, w \rangle} \right) \right]$

. מציאת פתרון ע"י לקיחת מינוס ושימוש באיטרציות ניוטון-רפסון. -סיתוך (cutoff): בגלל שאנו מעוניינים בתוצאה מתוך $\mathcal Y$ נבצע $\hat{y} = \begin{cases} 1 & h_w(x) \ge \alpha \\ 0 & h_w(x) < \alpha \end{cases}$

(Receiver Operating Characteristic עקומת ROC) ROC מתאימה בין TPR לבין FPR. אנו משתמשים במדד AUC (ר"ת . נתון w עבור α נתון (Area Under the Curve

צורה נוספת ל-Logistic Regression: (זו פונקציה קמורה!!) $\arg \min_{w} \frac{1}{m} \sum_{i=1}^{m} \log \left(1 + e^{-y_i \langle w, x_i \rangle}\right)$

Nearest Neighbors

k ועבור, $ho(x_i,z) = \sum_{j=1}^d w_j \cdot \left(x_j - y_j
ight)^2$ נגדיר מרחק ממושקל השכנים הקרובים נבחר:

$$\arg\max_{y\in\mathcal{Y}}\sum_{i=1}^k 1_{\{y_{\pi(i)}=y\}}$$

kd- עיבוד מקדים (למשל Brute Force), עיבוד מקדים (למשל tree) או חיפוש רנדומלי מהיר. (אין עקרון למידה).

$$\mathcal{H}^k_{\mathit{CT}} = \left\{ x \mapsto \sum_{j=1}^N c_j \cdot 1_{\{x_i \in B_j\}} \right\}$$

 $\mathbb{R}^d = \biguplus_j B_j$ כמות הקופסאות, k עומק מקסימליי. N-פונקציית התועלת: $P_y^{S}(B) = \frac{1}{n_S(B)} \sum_{x_i \in B} 1_{\{y_i = y\}}$, לפי עיקרון ה ERM מי שממקסם את ההסתברות זה כלל הרוב הפשוט.

> $g_i(t) = P_{\hat{y}\left(S(B_{>t})\right)}^S(B_{>t}) + P_{\hat{y}\left(S(B_{< t})\right)}^S(B_{< t})$ $\arg \max_{i,t} g_i(t)$

 $O(md \cdot 2^k) \stackrel{2^k \le m}{=} O(m^2d)$:זמן ריצה

 $L_{\mathcal{D},f}(h) = \Pr_{x \in \mathcal{D}}[h(x) \neq f(x)]$: Generalization Error

אלגוריתם למידה $\mathcal A$ בעל חסם:Probably Approximately $\delta>0$ ה מ-0 קטנה בהסתברות על (דיוק) $\varepsilon>0$ תחתון

(ביטחון). PAC איז היא וות $\mathcal H$ היא היפותזות $\mathcal H$ נאמר שמחלקת היפותזות אם קיימת פונקציה $\widetilde{m}_{\mathcal{H}} \colon (0,1)^2 o \mathbb{N}$ ואלגוריתם Learnable

למידה \mathcal{H} למידה \mathcal{H} בעל התכונה: $\mathcal{A}_m: (\mathcal{X} \times \mathcal{Y})^m \to \mathcal{H}$ למידה \bullet לכל $\varepsilon, \delta \in (0,1)$ \bullet לכל התפלגות \mathcal{D} מעל \mathcal{X} \bullet לכל פונקציה

(המקיימת את הנחת הריאלזביליות) $f \colon \mathcal{X} o \mathcal{Y}$ תיוג ${\mathcal D}$ כאשר נריץ את ${\mathcal A}$ על $\widetilde{m}_{\mathcal H}(arepsilon,\delta)$ דגימות i.i.d מעל מתוייגות ע"י f, מתקיים:

 $\Pr_{id_{n,m}}[L_{\mathcal{D},f}(\mathcal{A}(S)) > \varepsilon] < \delta$

 $\lim_{m \to \infty} \mathbb{E}_{\min_{s,m}[L_{D,f}(\mathcal{A}(\widetilde{S}))]} = 0$ עבור loss וסום, התכונה שקולה:

. המינימלי $\widetilde{m}_{\mathcal{H}}(arepsilon,\delta)$:Sample Complexity $arepsilon < rac{1}{2}$ אין ארוחות חינם: יהי \mathcal{X} מרחב דגימות אינסופי, ונקבע

קיימת $\delta > 0$ כך שלכל אלגוריתם למידה ${\cal A}$ וקורפוס אימון בגודל יכך ש: $f: \mathcal{X} \to \mathcal{Y}$ פונקצייה \mathcal{X} כך ש: $\mathcal{D}(x)$ מעל $\Pr_{\substack{i \not l \\ U_{\mathcal{D},f}}} \left[L_{\mathcal{D},f}(\mathcal{A}(S)) > \varepsilon \right] \ge \delta$

 $\exists h \in \mathcal{H} \ \Pr[h(x) = f(x)] = 1$:Realizability Assumption

תהי נסמן: ער היפותזות. נסמן: VC-Dimension $:\mathcal{H}_{c}=\{h\upharpoonright_{c}:h\in\mathcal{H}\}$ $VCdim(\mathcal{H}) = \max\{|C| : C \subset \mathcal{X} \text{ and } |\mathcal{H}_C| = 2^{|C|}\}$

הכללת המושגים:

עבור Approximately Correct • $\mathcal{X} imes \mathcal{Y}$ עבור • \mathcal{D} הבללת • $L_{\mathcal{D}}(h) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \varepsilon$ הבללת $h \in \mathcal{H}$ $L_{\mathcal{D}}(h) = \mathop{\mathbb{E}}_{(x,y) \sim \mathcal{D}}^{n} [\ell(h(x),y)]$ השגיאה ע"י

נאמר שמחלקת היפותזות ${\mathcal H}$ היא :Agnostic-PAC Learnable $\ell \colon \mathcal{H} \times (\mathcal{X} \times \mathcal{Y}) o [0, \infty)$ ביחס ל-Agnostic-PAC Learnable אם קיימת פונקציה $\widetilde{m}_{\mathcal{H}} \colon (0,1)^2 o \mathbb{N}$ ואלגוריתם למידה :בעל התכונה \mathcal{A}_m : $(\mathcal{X} \times \mathcal{Y})^m \to \mathcal{H}$

 $\mathcal{X} \times \mathcal{Y}$ לכל התפלגות \mathcal{D} מעל \bullet $\varepsilon, \delta \in (0,1)$ לכל \bullet כאשר נריץ את \mathcal{D} \star $m \geq \widetilde{m}_{\mathcal{H}}(\varepsilon, \delta)$ מעל \mathcal{A} אוו מעל מתקיים:

 $\Pr_{s^{lid} \cap m} \left[L_{\mathcal{D},f}(h_S) \le \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon \right] \ge 1 - \delta$

המשפט היסודי של הלמידה הסטטיסטית: $d = VCdim(\mathcal{H})$ מחלקת היפותזות ונסמן $\mathcal{H} \subset \{\pm 1\}^{\mathcal{X}}$ d < אם"ם אם Agnostic-PAC אם"ם אם PAC אז ${\mathcal H}$ למידה

PAC למידה \mathcal{H} -עך ש- \mathcal{H} למידה $\mathcal{C}_1,\mathcal{C}_2$ כך ש-א עם סיבוכיות דגימות:

 $C_1 \frac{d + \log\left(\frac{1}{\delta}\right)}{2} \le m_{\mathcal{H}}(\varepsilon, \delta) \le C_2$ $d\log\left(\frac{1}{\varepsilon}\right) + \log\left(\frac{1}{\delta}\right)$ למידה \mathcal{H} - פר כך $\mathcal{C}_1,\mathcal{C}_2$ למידה אוניברסליים קבועים אוניברסליים פ אם סיבוכיות דגימות: Agnostic-PAC

 $\frac{d + \log\left(\frac{1}{\delta}\right)}{2} \le m_{\mathcal{H}}(\varepsilon, \delta) \le C_2 \frac{\omega}{2}$ $d + \log\left(\frac{1}{\delta}\right)$ ε^2 • את החסם העליון ניתן להשיג ע"י לומד *ERM*

:מיימת S המקיימת, קבוצת דגימות S המקיימת, קבוצת דגימות $\forall h \in \mathcal{H} |L_S(h) - L_D(h)| < \varepsilon$

 $h_S \in \mathcal{A}, \mathcal{H}, \mathcal{D}$ למה f: תהי f קבוצת דגימות f-מייצגת עבור f $L_{\mathcal{D}}(h_{\mathcal{S}}) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon$ אז $ERM_{\mathcal{H}}(\mathcal{S})$

התכנסות אחידה: למחלקת היפותזות ${\mathcal H}$ יש תכונת התכנסות ולכל $arepsilon,\delta\in(0,1)$ כך שלכל $m^{\mathit{UC}}_{\mathcal{H}}:(0,1)^2 o\mathbb{N}$ ולכל $arepsilon,\delta\in(0,1)^2$ מתקיים: מעל \mathcal{D} מעל מעל מעל מעל - representative] ≥ 1 - δ

טענה: אם \mathcal{H} בעלת תכונת התכנסות אחידה עם פונקציה עם סיבוכיות Agnostic-PAC עם למידה \mathcal{H} אז \mathcal{H} אז $m_{\mathcal{H}}^{\mathit{UC}} \colon (0,1)^2 o \mathbb{N}$ $.m_{\mathcal{H}}(\varepsilon,\delta) \leq m_{\mathcal{H}}^{UC}\left(\frac{\varepsilon}{2},\delta\right)$ דגימות

הוכחת המשפט (המקרה הכללי)

יבורם: $b, \beta > 0 \;, m_0 \in \mathbb{N}$ עבורם אם קיימים לונומיאלית ב- $\forall m > m_0 \ \tau_{\mathcal{H}}(m) \coloneqq \max_{\substack{C \subset X \\ |C|}} |\mathcal{H}_C| \le b \cdot m^\beta$

 $au_{\mathcal{H}}(m) \leq \left(rac{em}{d}
ight)^d$ אז $m > VCdim(\mathcal{H})$ למה 3: אם

אז $|\mathcal{H}_c|$ גדל פולינומיאלית ב- $|\mathcal{C}|$, אז $|\mathcal{H}_c|$ מקיים את

תכונת ההתכנסות האחידה. $F_m^{\mathcal{D}}(S) \coloneqq \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$ נגדיר: צ"ל: להראות: $\Pr_{S_{i \to \mathcal{D}}^{i \to n}}^{n \in \mathcal{R}}[F_m^{\mathcal{D}}(S) > \varepsilon] < \delta$ לפי מרקוב מספיק להראות:

 $\Pr_{S \overset{iid}{\sim} \mathcal{D}^m} [F_m^{\mathcal{D}}(S) > \varepsilon] \leq \frac{\overset{iid}{\mathbb{E}} [F_m^{\mathcal{D}}(S)]}{\varepsilon^{iid} \mathcal{D}^m}$

 $\underset{S\overset{iid}{\sim}\mathcal{D}^m}{\mathbb{E}}[F^{\mathcal{D}}_m(S)] \leq O\left(\frac{\sqrt{\log(\tau_{\mathcal{H}}(2m))}}{\sqrt{2m}}\right) + o(m)$

<u>הוכחת המשפט (המקרה הסופי)</u> $\Pr_{S \stackrel{iid}{\sim} \mathcal{D}^m} [F_m^{\mathcal{D}}(S) > \varepsilon]$ def = $\Pr_{\substack{i \in \mathcal{D}_{D}^{m}}} [\exists h \in \mathcal{H} | L_{\mathcal{D}}(h) - \overline{L_{\mathcal{S}}(h)} | > \varepsilon$ $\sum_{i \in \mathcal{H}} \Pr_{S \stackrel{iid}{\sim} \mathcal{D}^m} [|L_{\mathcal{D}}(h) - L_{S}(h)| > \varepsilon]$ $|\mathcal{H}| \max_{h \in \mathcal{H}} \Pr_{S_{\underline{i}} \underline{i} \underline{d}_{\mathcal{D}} m} [|L_{\mathcal{D}}(h) - L_{S}(h)| >$

 $m \geq \frac{\log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2\varepsilon^2}$ אם ניקח $m \geq \frac{\log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2\varepsilon^2}$

 $\frac{$ מחלקות סופיות חוק $h\in\mathcal{H}$ המקיים $y_i=h(x_i)$ עבור $h\in\mathcal{H}$ המקיים חוק PAC משפט: כל מחלקת היפותזות \mathcal{H} מגודל סופי היא למידה ע"י שימוש בעקרון ה-ERM ובעלת סיבוכיות דגימה:

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \frac{\log(|\mathcal{H}|/\delta)}{\epsilon}$$

Number of examples (m)

Confidence (δ) Accuracy (ε)

 $\mathit{NVCdim}(\mathcal{H}) \leq \log_2 |\mathcal{H}|$, טענה: עבור \mathcal{H} סופית,

 σ^2 וקורלציה אונות X_1,\dots,X_T וקורלציה בהינתן מ"מ היא: $ar{X} = rac{1}{T}\sum_{t=1}^T X_t$ השונות של $.corr(X_i, X_j) =
ho$

:Meta-Algorithms בחירה ע"י רוב $h(x) = sign(\sum_t h_t(x))$ קלסיפיקציה •

 $h(x) = \frac{1}{x} \sum_t h_t(x)$ רגרסיה •

לתהליך ה-bootstrap שנוצר מ-S.

חסרונות: • צריך לאמׄן T מודלים. • בשביל חיזוי צריך לשמור

לא להשתמש כאשר הצלחת הלומד קטנה מחצי! The Bootstrap

אנו יוצרים B קבוצות דגימות S אנו יוצרים בהינתן קבוצת דגימות אנו יוצרים מתוך S במשך m פעמים: $S^b = \{(x_i^b, y_i^b)\}_{i=1}^m$

T מודלים. ● יותר קשה לפרש מדוע האלגוריתם חזה משהו.

מעל S י"י התפלגות $\widehat{\mathcal{D}}_{\mathcal{S}}$ התפלגות :Empirical Distribution $: \mathcal{C} \subset \mathcal{X} \times \mathcal{Y}$ על תת-קבוצה $\mathcal{X} \times \mathcal{Y}$

 $\widehat{\mathcal{D}}_{S}\big((X,Y)=(x,y)\big)=\begin{cases}\frac{1}{m} & (x,y)\in S\\ 0 & otherwise\end{cases}, \ \widehat{\mathcal{D}}_{S}(C)=\frac{|C\cap S|}{m}$ 0 otherwise ענה: מ"מ $X_1, ..., X_T$ ב"ת עם התפלגות F כלשהי מעל \mathbb{R} . יהי X_1 ערך דגימה של X_1 , אז CDF היא X_1 x_i בכל בקפיצות של מ-0 ל-1 בקפיצות של פונקציית מדרגות העולה מ-0 ב טענה: עבור m גדול דיו, בערך 37% מהנקודות נותרות מחוץ

(הוכדת ה-Bias ה-Variance נותר זהה) Bagging עבור אלגוריתם למידה $\mathcal A$ בסיסי בוחרים T ויוצרים מדגמים h_{S^1}, \dots, h_{S^T} כל אחד בגודל m. מאמנים T היפותזות S^1, \dots, S^T חיזוי דגימה חדשה: $h_{bag}(x) = sign(\sum_{t=1}^{T} h_{S^t}(x))$

De-correlation: מכיווו שניתו להוריד את השונות רק עד מידה מסויימת, אנו מבצעים תהליך דה-קורולציה. למשל בעצים נבחר . פרמטר שרק מתוכם ניתן לפצל k ובכל עת נבחר k

(הורדת ה-Bias, ה-Variance עולה לאט) Boosting ע"י עדכון מרחב המדגם S^t (הוספת משקולות) כך שדגימות בהן U^t טעינו יקבלו משקל גדול יותר (יוצר התפלגות D^t). [או שניתן לדגום מחדש מתוך התפלגות \mathcal{D}^t , תמיד זמין אך יכול ליצור פילויות דגימה]

Adaboost Algorithm:

Initialize:
$$D^1 = \left(\frac{1}{m}, \dots, \frac{1}{m}\right)$$

Loop: $h_t = \mathcal{A}(D^t, S)$

Update:

 $\bullet \ \sum_{j=1}^{m} D_j^{t+1} \cdot 1_{\left[h_t(x_j) \neq y_j\right]} = \frac{1}{2} \Rightarrow w_t = \frac{1}{2} \log \left(\frac{1}{\varepsilon_t} - 1\right)$ $(\varepsilon_t = \sum_{j=1}^m D_j^t \cdot \mathbf{1}_{[h_t(x_j) \neq y]}^t$ י (פיכון אמפירי ממושקל $\mathbf{1}_{[h_t(x_j) \neq y]}^t$ י מושקל $\mathbf{0}_t^{t+1} = \frac{D_t^t \cdot e^{-w^t y_t h_t(x_t)}}{\mathbf{1}_{[t+1]}^t}$

 $\sum_{i=1}^{m} D_{i}^{t} \cdot e^{-w^{t}y_{j}h_{t}(x_{j})}$

זיזוי דגימה חדשה:

$$h_{boost}(x) = sign(\sum_{t=1}^{T} w_t h_t(x))$$

אלגוריתם למידה $\mathcal A$ הוא לומד- γ -weak-learne מחלקת היפותזות \mathcal{H} אם קיימת פונקציה \mathbb{N} כך $m_{\mathcal{H}}:(0,1) o \mathbb{N}$ שלכל (0,1) \mathcal{X} , לכל התפלגות \mathcal{D} מעל מרחב מדגם \mathcal{X} , ולכל מתקיימת מתקיית תיוג $f \colon \mathcal{X} \to \{\pm\}$, אם הנחת הריאלזביליות מתקיימת $m\geq m_{\mathcal{H}}(\delta)$ על \mathcal{A} על נריץ את געשר נריץ אז כאשר ($\mathcal{H},\mathcal{D},f$ -ל) ביחס ל $h_S=\mathcal{A}(S)$ מתוך אימות $\mathcal{D}\times f$ האלגוריתם יחזיר היפותזה

$$\Pr_{\substack{S \stackrel{iid}{\sim} \mathcal{D}^m}} \left[L_{\mathcal{D},f}(h_S) \le \frac{1}{2} - \gamma \right] \ge 1 - \delta$$

 $S \sim \nu^{\prime\prime\prime}$ מחלקה אם קיים: γ -weak-learnable

 $VCdim(\mathcal{H}) < \infty \Leftarrow$ משפט: למידות חלשה

משפט: תהי S קבוצת דגימות. נניח שבכל איטרציה של עבודה (h_t הלומד הבסיסי מחזיר כלל חיזוי (היפותזה Adaboost

הסיכון האמפירי הממושקל מקיים $\varepsilon_t = \sum_{j=1}^m D_j^t \cdot 1_{\left[h_t\left(x_j\right) \neq y_j\right]} \leq \frac{1}{2} - \gamma$ אז כלל הסיכון האמפירי (הלא ממושקל) של כלל ההחלטה של

ומקיים: Adaboost
$$L_S(h_{boost})\equiv rac{1}{m}\sum_{j=1}^m 1_{\left[h_{boost}(x_j)\neq y_j\right]} \leq e^{-2\gamma^2 T}$$

	Bagging	Boosting
מקביליות	במקביל	טורי
מבנה הנתונים	Bootstrap training samples	ממושקל או Bootstrap קבוצה S המקורית עם כלל ERM ממושקל
De-correlation	מומלץ	אין צורך
Overfitting	לא קורה	סיכון
איזה מודל להשתמש	עצים עמוקים	עצים רדודים
השפעה	מפחית Variance	מפחית Bias

Regression & Regularization

מהנוסחה לעיל הוא תנאי הרגולציה. \mathcal{R} :Regularization Term אין התייחסות כלל למורכבות ההיפותזה $\lambda=0$ שנבחרה. ה-Bias נמוך אך ה-Variance גבוה.

עבור $\infty o \lambda$ נעדיף היפותזה כמה שיותר פשוטה ללא קשר λ . לפונקציית המטרה.

Regression Trees

עבור עץ רגרסיה |T| את גדיר את $h(x) = \sum_{j=1}^N c_j \cdot 1_{[x \in B_j]}$ להיות מספר העלים, ואת T_0 להיות העץ המתקבל מתהליך הגידול: $\min_{T\subseteq T_0} L_S(T) + \lambda \cdot |T|$

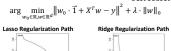
Modern Regression

Best Subset •

Ridge Regression • $\min_{\mathbf{x} \in \mathbb{R}, w \in \mathbb{R}^d} \| w_0 \cdot \vec{1} + X^T w - y \|^2 + \lambda \cdot \| w \|_2^2$

 $\widehat{w}_{\lambda}=$ יש לפתור את המערכת $Xy=(XX^T+\lambda I)w$ פתרון יהיה של לפתור את המערכת $\Sigma^{\lambda}=diag\left(rac{\sigma_i}{\sigma_i^2+\lambda}
ight)$ כאשר ער משר $U\Sigma^{\lambda}V^Ty$ • Lasso – בעל תכונת דלילות.

$\arg\min_{w_0\in\mathbb{R},w\in\mathbb{R}^d}\left\|w_0\cdot\vec{1}+X^Tw-y\right\|^2+\lambda\cdot\|w\|_1$



: טענה: עבור אורתוגונלי) $XX^T=I$ ומטריצה $w\in\mathbb{R}^d$ עבור $\hat{w}_\lambda^{ridge}=\frac{1}{1+\lambda}\hat{w}^{LS}$

 $\mathrm{s.t}\,\eta_{\lambda}^{soft}(x) = \begin{cases} x-\lambda & x \geq \\ 0 & -\lambda < x < \\ x+\lambda & -\lambda \geq x \end{cases}$ $\bullet \ \widehat{w}_{\lambda}^{lasso} = \eta_{\lambda}^{soft}(\widehat{w}^{LS})$ • $\widehat{w}_{\lambda}^{subset} = \eta_{\sqrt{\lambda}}^{hard}(\widehat{w}^{LS})$

Model Selection

Regularized Logistic Regression $\arg\min_{w_0,w} \left[\sum \left[\log \left(1 + e^{w_0 + \langle x_i, w \rangle} \right) - y_i(w_0 + \langle x_i, w \rangle) \right] + \lambda \cdot \|w\|_1 \right]$

-ה עבור את המודל $i=1,\ldots,k$ עבור: k-fold Cross Validation של loss- נחשב את ה-i. נחשב את ה-מלבד החלק של . המודל ה-i ע"פ החלק ה-i . נחזיר את הממוצע וסטיית התקן של התוצאות.

$$\sigma = \sqrt{\frac{1}{n-1}\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

. נדגום B קבוצות S^b , ונאמן את האלגוריתם עליהם Bootstrap ונמצע $T^b = S \setminus S^b$ ונמצע בעדי לבחון את הביצועים נשתמש – כמקודם.

:טעויות נפוצות

- 1. Over-estimating generalization error כאשר מאמנים PAC על $\frac{k-1}{k} \cdot m$ כמות לא מספקת של דגימות, אז לפי .נעריך בחסר את השגיאה
- כאשר אנו Under-estimating generalization error .2 מעריכים מודל ע"פ נתונים עליהם אומן (מצב של Overfitting), או כאשר על הדאטה בוצע עיסוי המתאים למודל.

$$L_{\mathcal{D}}(h_{\mathcal{S}}) = \underbrace{L_{\mathcal{D}}(h_{\mathcal{S}}) - L_{\mathcal{V}}(h_{\mathcal{S}})}_{} + \underbrace{L_{\mathcal{V}}(h_{\mathcal{S}}) - L_{\mathcal{S}}(h_{\mathcal{S}})}_{} + \underbrace{L_{\mathcal{S}}(h_{\mathcal{S}})}_{} + \underbrace{L_{\mathcal{S}}(h_{\mathcal{S}})}_{}$$

A. שגיאת ההכללה. ניתנת לחסימה תחת ההנחה שפונקציית ה-.Overfitting אבור B גדול ו- \mathcal{C} קטן ככל הנראה B. עבור B. .Underfitting עבור C גדול ככל הנראה C

$$\delta \in (0,1)$$
 ולכל ולכל $h \in \mathcal{H}$ טענה: לכל

$$\Pr\left[|L_{\mathcal{V}}(h) - L_{\mathcal{D}}(h)| \le \sqrt{\frac{\ln(2/\delta)}{m_{\mathcal{V}}}}\right] \ge 1 - \delta$$

עבור $\mathcal{H}_1 \subseteq \cdots \subseteq \mathcal{H}_k$ סופיים. $:h^* \in \mathit{ERM}_{\mathcal{H}_k}(\mathcal{S}_{all})$ בבור יבור יבו

$$\Pr\left[L_{\mathcal{D}}(h^*) \leq \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{21 \left(\frac{2|\mathcal{H}_k|}{\delta}\right)}{m}}\right] \geq 1 - \delta$$

 $h^* \in \arg\min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) \subseteq \mathcal{H}_j \subseteq \cdots \subseteq \mathcal{H}_k$ עבור :Model Selection

$$\Pr\left[L_{\mathcal{D}}(h^*) \leq \min_{h \in \mathcal{H}_k} L_{\mathcal{D}}(h) + \sqrt{\frac{2}{am} \ln\left(\frac{4k}{\delta}\right)} + \sqrt{\frac{2}{(1-\alpha)m} \ln\left(\frac{4|\mathcal{H}_j|}{\delta}\right)}\right] \geq 1 - \delta$$

בחירת פיטצ'רים 1. מתאם פירסון

$$R = \frac{\sum_{i} (x_{i} - \bar{x})(y_{i} - \bar{y})}{\sqrt{\sum_{i} (x_{i} - \bar{x})\sum_{i} (y_{i} - \bar{y})}}$$
DDDD A - Forward-Stopwice Sole

- 2. Forward-Stepwise Selection אלגוריתם חמדן שמתחיל k ומוסיף פיטצ'רים עד לכמות (intercept) מהחותך באלגוריתם המתחיל מכל - Backward-Stepwise Selection . 3
- הפרמטרים ומוריד כל פעם פרמטר עם המתאם הקטן 4. Regularization - שיטה כללית בה מגבילים את כמות
- בנוסף ניתן לבצע שינוי של המידע, למשל מרכוז שלו סביב

הראשית, שינוי הטווח (Scaling), ביצוע Standardization ועוד..

<u>דוגמאות לשימוש:</u> 1. חשיפה של מבנה במימד נמוך. 2. קלסיפיקציה. 3. זיהוי אנומליות.

 $\arg\min_{U \in \mathbb{R}^{d \times k}, W \in \mathbb{R}^{k \times d}} \sum_{i=1}^{n} ||x_i - UWx_i||^2$ A משפט: תהי $u_1, ..., u_n$ ויהיו, $A = \sum_{i=1}^m x_i x_i^T$ הו"ע של

פתרון לבעיית $U=W^T=[u_1\cdots u_k]$ פתרון לבעיית) נקראת מטריצת הטלה, היא סימטרית והיא $P = \sum_{i=1}^k u_i u_i^T ullet$

- ייין ביותר ביותר ל-x בתת-מרחב:
- $\|x-v\|_2 \geq \|x-Px\|_2 \ \forall v \in span\{u_1,\dots,u_k\}$ בלבד. y בלבד בו מטילים את הערך בלבד.
- ע"י הזזה של הנקודות $ar{x}_i = x_i ar{x}$ אנו מקבלים ב-4 את $A = rac{1}{m-1} \sum_i (x_i - ar{x}) (x_i - ar{x})^T$ מטריצה השונות:

 $oldsymbol{x}_m$ **הגדרות:** תהי מטריצת השונות של דגימות אימון כמוגדר לעיל, ויהיו $u_1,...,u_d$ הו"ע של מטריצת השונות $\lambda_1\geq\cdots\geq\lambda_d\geq 0$.

- x_1,\dots,x_m נקרא הערך המוביל ה-i של נקרא λ_i נקרא •
- $x_1,...,x_m$ של i-ה הוקטור ול נקרא הוקטור במרחב היל ניתן לייצוג ullet הורדת מימד תהיה U^Tx_i לאחר מכן כל וקטור במרחב ניתן לייצוג $ar{x} + \sum \alpha_i u_i$ ע"י הקירוב

PCA Algorithm:

Input: $X \in \mathbb{R}^{m \times d}$ Eval:

o $A = X^T X$

0 R - X R $0 u_1, ..., u_k$ eigenvectors of A• else $(m \le d)$:

o $B = XX^T$ o $v_{\text{\tiny 1}},\dots,v_{k}$ eigenvectors of B

o denote $u_i = \frac{1}{\|X^T v_i\|} X^T v_i$

• return u_1, \dots, u_k

. בשביל לבחור את k ניתן לייצר Scree Plot (גרף בו מציגים את הע"ע בסדר יורד).

נגדיר חלוקה של המידע $x_1,...,x_m$ ל-k מחלקות $\bigcup_{i} C_i$. נגדיר ע"י: ערחק ע"י (d מרחק מרחק ע"י:

יית המחיר (ע"פ פונקציית מרחק ע"י.
$$G(C_1, ..., C_k) = \min_{\mu_1, ..., \mu_k} \sum_{j=1}^k \sum_{x \in C_j} d(x, \mu_j)^2$$

. arg $\min_{u} \sum_{x \in C_j} d(x, \mu_j)^2$ הערך (C_j של :Centroid

centroid- טענה: אם $\mathcal{X} = \mathbb{R}^d$ ו- \mathcal{X} היא הנורמה האוקלידית, אז $.\mu_j = \frac{1}{|c_j|} \sum_{x \in c_j} x$ הוא הממוצע

k-means Algorithm:

Input: $x_1, ... x_m$ and $k \in \mathbb{N}$ Step 0: choose initial $\mu_1, ..., \mu_k$ Until convergence:

- ullet set \mathcal{C}_{j} to be the points x_{i} closer to μ_i than to any other centroid.
- update μ_j to centroid of C_j :

$$u_j = \frac{1}{|C_j|} \sum_{x \in C_j} x$$

..., μ_{ν} נקראת Voronoi cells החלוקה המושרית מ $\mu_{\nu}, \dots, \mu_{\nu}$

תמיד השונות בכל \mathcal{C}_j יורדת בכל איטרציה. האלגוריתם תמיד מתכנס. האלגוריתם מתכנס למינימום מקומי (בהתאם לנקודות ההתחלה).

 $f \colon \mathcal{C} o \mathbb{R}$ פונקציה קמורה. פונקציה \mathcal{C} קבוצה קמורה. פונקציה נקראת קמורה אם:

$$\forall u, v \in \mathcal{C}, \alpha \in [0,1] \ f(\alpha v + (1-\alpha)u)$$

$$\leq \alpha f(u) + (1-\alpha)f(v)$$

 $\bigcap_{lpha} A_lpha$ חיתוך של קבוצות קמורות הוא קמור. משפט:

 $epi(f) = \{(x, t) | f(x) \le t\}$ אפיגרף: משפט: פונקציה f קמורה אם"ם epi(f) קמורה אם"ם $f(y) \geq f(x) + \nabla f(x)^T (y-x)$ מתקיים $x,y \in dom(f)$

 $\nabla^2 f(x) \geqslant 0$ מסדר ראשון) אם"ם $\nabla^2 f(x) > 0$ (תנאי מסדר שני).

תכונות השומרות על קמירות:

- אז , $\alpha_1,\dots,\alpha_n\geq 0$, קמורות, f_1,\dots,f_n אבור סכום חיובי . קמורה $f = \sum_i \alpha_i f_i$
- קמורה על g(x) = f(Ax + b) קמורה על $.\{x|Ax+b\in dom(f)\}$
- x- מהסימום אם לכל $v \in A$ מתקיים f(x, v) קמורה ב-3. . אז $g(x) = \sup f(x, y)$ אז
- S-ו ($\nabla^2 f \geqslant 0$) אורה במשותף (f(x,y) ו-קמורה. (לא בקורס) קמורה, אז $g(x) = \inf_{y \in S} f(x,y)$ קמורה

תכונות - מינימום משפט: אם f קמורה, כל מינימום מקומי הוא גלובלי. משפט: תה f קמורה ודפרנציבילית ו- $w\in dom(f)$, אז: $\forall u \ f(u) \ge f(w) + \langle \nabla f(w), u - w \rangle$

תכונות – סאב-גרדיאנט

 \overline{v} אם: \overline{v} אם: אם אב-גרדיאנט של $orall u f(u) \geq f(w) + \langle v, u - w \rangle$ סימון: $\partial f(w)$ זוהי קבוצת כל סאב-הגרדיאנטים של $\partial f(w)$ מתקיים $\partial f(w)$ ממה: $\partial f(w)$ ממקיים $\partial f(w)$

 $.\partial f(x) = \{\nabla f(x)\}$ אם f דפרנציבילית ב-x אז $\partial f(w) \neq \emptyset$ טענה: נניח $f(v)=\max f_i(v)$ קמורות, ותהי $f_i\colon V o \mathbb{R}$ כמו $\partial f_j(u) \subseteq \partial f(u)$ אז $j \in \arg\max f_i(u)$ נסמן $u \in V$ כן עבור $ec{0}\in\partial f(u)$ עבורה $w\in V$ קמורה. תהי $f\colon V o\mathbb{R}$ עבורה f אז w מינימום גלובלי של w

תכונות - ליפשיץ $f\colon \mathcal{C} o \mathbb{R}$ נקראת ho-ליפשיץ אם: $\forall w_1, w_2 \in \mathcal{C} \ |f(w_1) - f(w_2)| \leq \rho \|w_1 - w_2\|$ למה: אם f המורה. אז f היא ρ -לפישיץ אם"ם הנורמה של כל ho סאב-גרדיאנט של f הוא לכל היותר

אופטימיזציה קמורה

 $\min_{\mathbf{x}} f_0(\mathbf{x})$ בעיית אופטימיזציה: $s.t \quad f_i(x) \leq b_i \quad \forall i=1,...,n$

. בעיית אופטימיזציה f_i קמורות אופטימיזציה עם בעיית אופטימיזציה עם . לינאריות f_i עם בעיית אופטימיזציה עם לינאריות בעיית תכנון לינארי

תקרא $\mathcal{X} \times \mathcal{Y}$ מעל \mathcal{H}, ℓ מעל בעיית למידה מחורה: בעיית למידה מחלקת ההיפותזות \mathcal{H} היא קבוצה קמורה, ולכל h-ב הפונקציה $\ell(h(x),y)$ הפונקציה $(x,y) \in \mathcal{X} \times \mathcal{Y}$. $ERM_{\mathcal{H}}$ דוגמה:

.PAC-טענה: לא כל הבעיות הקמורות מעל \mathbb{R}^d הן למידות יין היים אינים אינים אינים \mathcal{H}^d , כך ש \mathcal{H}^d חסומה ו ℓ^+ ליפשיץ, היא

- $\forall w \in \mathcal{H} \ \|w\| \leq B$ חסומה קיים B עבורו $\mathcal{H} ullet$
- ho והיא h-ם הפונקציה $\ell(h(x),y)$ קמורה ב- ℓ ליפנעיע $\epsilon,\delta,B,
 ho$ במקרה זה סיבוכיות הדגימה תלויה רק ב-

Sub-gradient Descent Algorithm:

initialize: $w^{(1)} = 0$ for $t=1,\ldots,T$:

• Choose $v_t \in \partial f(w^{(t)})$

 $\bullet \ w^{(t+1)} \leftarrow w^{(t)} - \eta \cdot v_t$

return $\overline{w} = \frac{1}{T} \sum_{t=1}^{T} w^{(t)}$

 $w^* \in \mathcal{P}$ ויהי ויהי פונקציה קמורה, ρ -ליפשיץ, ויהי T על אם נריץ את האלגוריתם Sub-GD אם נריץ את במשך. arg $\min f(w)$:צעדים עם \overline{w} אז הוקטור \overline{w} אז הוקטור מקיים:

$$f(\overline{w}) \le f(w^*) + \frac{\|w^*\|\rho}{\sqrt{T}}$$

 $w^*\in {}^{v^*}$ מסקנה: תהי f פונקציה קמורה, ho-ליפשיץ, ויהי על Sub-GD אם נריץ את האלגוריתם arepsilon>0 על מריתם arepsilon>0 אם נריץ את האלגוריתם במשך $rac{\|w^*\|^2
ho^2}{arepsilon^2}$ צעדים עם $\frac{\|w^*\|^2
ho^2}{
ho \sqrt{T}}$, אז הוקטור \overline{w} המוחזר f

 $f(\overline{w}) \leq f(w^*) + \varepsilon$ איטרציות בכדי Sub-GD צריך איטרציות בכדי Sub-GD מסקנה: האלגוריתם

Stochastic Gradient Descent Algorithm:

initialize: $w^{(1)} = 0$

- for t = 1, ..., T:
 - Choose $(x, y) \sim \mathcal{D}$
- Choose $v_t \in \partial \ell \left(w^{(t)}, (x, y) \right)$ $\bullet \ w^{(t+1)} \leftarrow w^{(t)} - \eta \stackrel{\backprime}{\cdot} v_t$
- return $\overline{w} = \frac{1}{T} \sum_{t=1}^{T} w^{(t)}$

למה: תהי בעיית למידה קמורה-ליפשיץ-חסומה עם פרמטרים למזעור SGD אז לכל $\epsilon>0$, אם נריץ את האלגוריתם. ho,Bעם מספר איטרציות $T \geq \frac{B^2 \rho^2}{\varepsilon^2}$ ועם $T \geq \frac{B^2 \rho^2}{\varepsilon^2}$ עם מספר איטרציות עם מספר איטרציות של SGD מקיים:

 $\mathbb{E}[L_{\mathcal{D}}(\overline{w})] \leq \min_{\sigma \in L} L_{\mathcal{D}}(w) + \varepsilon$

. אלגוריתם ${\mathcal A}$ הלומד בעזרת SGD טענה: אלגוריתם

טבלת VCdim: VCdim Axis aligned rectangles Homogeneous Halfspaces d. d + 1Non- Homogeneous Halfspaces Polynomial Thresholds (Pol.,) $\leq {d+r-1 \choose d}$ k-interval 2kConjuctions

תזכורת: בשביל $VCdim(\mathcal{H})=d$ צריך להראות שקיימת מחלקה מגודל d+1 לא ניתן 'נתץ.

ההבדלים בשגיאות:

<mark>כדאי לזכור</mark>

- $\bullet \ {\bf Generalization} \ {\bf error} : L_{\mathcal D}(h) \\$
- ullet Training error / empirical error: $L_S(h)$
- Approximation error: $\min L_{\mathcal{D}}(h)$
- Estimated error: $L_{\mathcal{D}}(h) \min_{i} L_{\mathcal{D}}(h')$

כלומר את מספר (כלומר את $\overset{..}{T}$ ב-Adaboost כלומר את מספר הלומדים החלשים) אנו נשפר את ה-Approximation error (Estimated error-).

 $h_S = ,h^* \in \arg\min_h L_{\mathcal{D}}(h)$ יהי: Bias-Variance Tradeoff

$$L_{\mathcal{D}}(h_{\mathcal{S}}) = \underbrace{L_{\mathcal{D}}(h^*)}_{\varepsilon_{annyroximation}} + \underbrace{L_{\mathcal{D}}(h_{\mathcal{S}}) - L_{\mathcal{D}}(h^*)}_{\varepsilon_{extimation}}$$

	KNN		
	Low Bias	Small $\stackrel{k}{\rightarrow}$ Big	High Bias
	Low Var	Big ← Small	High Var
Ì			

Decision Tree		
Low Bias	Deep ^{height} Shallow	High Bias
Low Var	Shallow \xrightarrow{height} Deep	High Var

Linear Regression			
Low Bias	(d~m)	$(d \ll m)$	High Bias
Low Var	$(d \ll m)$	(d~m)	High Var

Regularization		
Low Bias $0 \xrightarrow{\lambda} \infty$ High Bias		High Bias
Low Var	$\infty \stackrel{\lambda}{\leftarrow} 0$	High Var

Bagging		
Bias	Dosen't change	Bias
Low Var	$Big(\infty) \stackrel{T}{\leftarrow} Small$	High Var

Adaboost		
Low Bias	$Big(\infty) \stackrel{T}{\leftarrow} Small$	High Bias
Low Var	Small $\stackrel{T}{\rightarrow}$ Big(∞)	Var goes up a bit

Clustering		
Low Bias Big ← Small High Bias		
Low Var	Small $\stackrel{k}{\rightarrow}$ Big	High Var

אלגברה לינארית

 $A \in \mathbb{R}^{m imes n}$ טענה: עבור מטריצה

$$\mathbb{R}^m = \mathcal{C}(A) \oplus \mathcal{N}(A^T), \qquad \mathbb{R}^n = \mathcal{C}(A^T) \oplus \mathcal{N}(A)$$

:משפט: עבור מטריצה $A \in \mathbb{R}^{m imes n}$ משפט: עבור מטריצה (rank(A) = m) מדרגה מלאה (A .1

 $det(A) \neq 0$.2 $im(A) = \mathbb{R}^m$.3

 $\ker(A) = \{\vec{0}\} .4$

 $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \Rightarrow A^{-1} = \frac{1}{\det(A)} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$ $(AA^T = A^TA = I)$ מטריצת סיבוב: מטריצה אורתוגונלית det(A) = 1 עבורה.

 $A = \begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix}$ $P_u(v) = \|v\|\cos(\Theta)\cdot rac{u}{\|u\|} = rac{\langle v,u
angle}{\|u\|^2}\cdot u$:הטלה

 $\hat{\mu} = \frac{1}{m} \sum_{i=1}^{m} x_i$ שערוך תוחלת: $\frac{1}{m-1}\sum_{i=1}^{m}(x_i-\hat{\mu})^2$ שערוך שונות:

שערוך מטריצת Covariance:

 $\hat{C} = \frac{1}{m-1} \sum_{i=1}^{m} (x_i - \bar{x})(x_i - \bar{x})^T$ מתאם פירסון: בין X,Y מ"מ: $\rho_{X,Y} = \frac{Cov(X,Y)}{Cov(X,Y)}$

 $\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$ יישוויון מרקוב: עבור X מ"מ אי-שוויון מרקוב: $\Pr[|X - \mathbb{E}[X]| \ge t] \le \frac{Var[X]}{t^2}$ אי-שוויון צ'בישב:

 $a_i \leq X_i \leq b_i$ ב"ת כך ש- אי-שוויון הופדינג: עבור X_1, \dots, X_n ב"ת כך אי-שוויון הופדינג $\Pr[|X - \mathbb{E}[X]| \ge t] \le 2e^{-\frac{2m^2t^2}{\sum_{i=1}^{m}(b_i - a_i)^2}}$

:סענה: עבור מטריצה A (קבועה), ווקטור מקרי z $Var[Az] = A \cdot Var[z] \cdot A^T$

$$J_{x}(f) = \begin{pmatrix} \frac{\sigma_{I1}}{\sigma_{x_{1}}} & \cdots & \frac{\sigma_{Im}}{\sigma_{x_{1}}} \\ \vdots & \ddots & \vdots \\ \frac{\sigma_{I}}{\sigma_{x_{n}}} & \cdots & \frac{\sigma_{Im}}{\sigma_{x_{n}}} \end{pmatrix} = \begin{pmatrix} 1 & \cdots & \nabla_{f_{n}} \\ \nabla f_{1} & \cdots & \nabla f_{n} \end{pmatrix}$$
יעקוביאן:

$$J_{x}(g \circ f) = J_{f(x)}(g) \cdot J_{x}(f)$$
 כלל השרשרת:

$$H(f) = \begin{pmatrix} \frac{\partial^2 f}{\partial^2 x_1} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_2} \end{pmatrix}$$

$$\nabla (||w||^2) = 2w, \quad J_w(A \cdot w) = A, \quad J_x(f) = (\nabla f)^T$$