

- השוני במודל האגנוסטי של PAC הוא שההתפלגות  $D$  היא על  $Z = X \times Y$ , לא חייבת להיות פונקציה  $f \in \mathcal{H}$ , והמטרה היא למזער עד כדי השגיאה המינימלית:  $L_D(A(S)) \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$ . נשים לב כי מתקיים  $D_x(y) = D(y|x) = D((x, y)|x)$ .
- בנוסף, במודל זה  $Y$  יכולה להיות קבוצה עם יותר איברים (קבוצה סופית בבעיית Multiclass או אינסופית בבעיית Regression). כתוצאה מכך, גם  $\ell(h, (x, y))$  יכולה להיות אחרת (ריבועי, ערך מוחלט, טבלת מחירים על  $Y \times Y$  וכו'). כך  $L_D(h) = \mathbb{E}[\ell(h, z)]$ .  $\mathcal{H}, Z, \ell$  הלומד יודע את  $\mathcal{H}, Z, \ell$ .
- מדגם נקרא  $\epsilon$ -מייצג אם לכל  $h \in \mathcal{H}$  מתקיים  $|L_S(h) - L_D(h)| \leq \epsilon$ . אם מדגם הוא  $\frac{\epsilon}{2}$  מייצג, אזי כל פלט של  $h_S \in \arg \min_{h \in \mathcal{H}} L_S(h)$  ERM  $\mathcal{H}(S) = h_S$  מקיים  $L_D(h_S) \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$ .
- למחלקה  $\mathcal{H}$  יש תכונה UC (Uniform Convergence) אם קיימות  $m_{\mathcal{H}}^{UC}(\epsilon, \delta)$  כך שלכל  $\epsilon, \delta$  והסתברות  $D$  מתקיים בהסתברות  $1 - \delta$  שמדגם בגודל  $m$  הוא  $\epsilon$  מייצג.
- במידה ולמחלקה  $\mathcal{H}$  יש את התכונה UC אזי היא למידת PAC אגנוסטי עם  $m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}(\epsilon/2, \delta)$  ולאגל ERM  $\mathcal{H}$  (כלומר התכונה UC מספיקה ללמידות).
- משפט: אם  $\mathcal{H}$  היא מחלקה סופית עם Loss חסום בקטע  $[0, 1]$  אזי  $\mathcal{H}$  היא למידת PAC אגנוסטי עם  $m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{2 \log(2|\mathcal{H}|/\delta)}{\epsilon^2} \right\rceil$  (כלומר יש לה את תכונת UC).
- שימוש במשפט הנ"ל וטריק הדיסקרטיזציה מאפשר ללמוד כל מחלקה שמאופיינת עם  $d$  מספרים שכל אחד מהם בייצוג של  $b$  ביטים בסיבוכיות מדגם  $m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{2db + 2 \log(2/\delta)}{\epsilon^2} \right\rceil$ .

## 2. PAC learning

- בלמידת Batch מקבלים קבוצת דוגמאות  $S \in (X \times Y)^m$ . הלומד מחזיר כלל תיוג  $h$  כרגיל. המטרה היא לצדוק על דוגמאות עתידיות, כלומר להחזיר  $h$  שדומה כמה שיותר ל  $f$  האמיתית.
- מגדירים טעות הכללה  $L_{D,f}(h) = \mathbb{P}[h(x) \neq f(x)] = D(\{x \in X : h(x) \neq f(x)\})$ . המטרה היא למזער את הפונקציה הזאת. (הדוגמאות הן i.i.d וכן מניחים Realizability (קיימת פונקציית תיוג אמיתית).
- מודל PAC - מנסים למזער את  $L_{D,f}$  עד כדי  $\epsilon$  ובסיכוי  $1 - \delta$  ע"פ בחירת המדגם (פרמטרים של המשתמש), כאשר  $D, f$  לא ידועים ללומד, וכן הוא יכול לבקש לפחות  $m(\epsilon, \delta)$  דוגמאות.
- משפט ה NFL אומר שאם  $\mathcal{H} = Y^X$  הבעיה בלתי פתירה, כלומר לכל אלגוריתם לומד יש  $D, f$  שבסיכוי גדול מ  $\delta$  הטעות עליהם תהיה יותר מ  $\epsilon$  (וכמובן קטנה מחצי - ניהוש בין שתי אופציות).
- למידת מחלקות סופיות באמצעות כלל למידה Consistent, נקרא גם ERM - מזעור שגיאה אמפירית (שגיאת אימון ולא הכללה).
- עבור מחלקה סופית אם  $m \geq \frac{\log(|\mathcal{H}|/\delta)}{\epsilon}$  אזי לכל  $D, f$  מתקיים בסיכוי של לפחות  $1 - \delta$  על בחירת  $S$  מגודל  $m$ :  $L_{D,f}(\text{ERM}_{\mathcal{H}}(S)) \leq \epsilon$  (הוכחה עם חסם האיחוד).
- מחלקה  $\mathcal{H}$  היא למידת PAC אם קיימת פונקציה  $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$  ולאגל כך שלכל  $\epsilon, \delta$  והסתברות  $D$  פונקציית תיוג  $f : X \rightarrow \{0, 1\}$ , כאשר מריצים את האגל על  $m_{\mathcal{H}}(\epsilon, \delta)$  דוגמאות i.i.d שנוצרו ע"י  $D$  ותיוג ע"י  $f$ , אזי האגל יחזיר היפותזה  $h$  כך שבסיכוי של לפחות  $1 - \delta$  מתקיים  $L_{D,f}(h) \leq \epsilon$ .  $m_{\mathcal{H}}$  (המינימלית) היא סיבוכיות המדגם של  $\mathcal{H}$ .
- עבור קבוצת דוגמאות  $\mathcal{H}, C \subset X$  מנתצת את  $C$  אם היא מכילה היפותזה לכל תיוג אפשרי של איברי  $C$ . מימד ה VC של מחלקה  $\mathcal{H}$  הוא גודל הקבוצה הכי גדולה שמנותצת ע"י  $\mathcal{H}$  (אם  $\mathcal{H}$  מנתצת קבוצה  $C$  כלשהי, היא לא מספקת עליה שום ידע מוקדם). למחלקות סופיות מתקיים  $\text{VCdim}(\mathcal{H}) \leq \log_2(|\mathcal{H}|)$ . המימד הוא בדר"כ מספר הפרמטרים שמאפיינים את המחלקה.
- אם  $\mathcal{H}$  היא מחלקת היפותזות של תיוג בינארי אזי קיימים  $C_1, C_2$  (מוחלטים) כך שסיבוכיות המדגם של למידת PAC של  $\mathcal{H}$  היא  $m_{\mathcal{H}}(\epsilon, \delta) \leq C_1 \frac{d + \log(1/\delta)}{\epsilon} + C_2 \frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}$ ,  $\text{VCdim}(\mathcal{H}) < \infty$  אם ורק אם.
- ERM עבור חצאי מרחב - אפשר עם אליפסואיד, או עם אלג' פרספטרון (מתחילת Boosting ממשקולות 0 ומעדכן את המשקולות עם  $y_i x_i$  כל עוד יש דוגמה שטועים עליה) אשר עושה לכל היותר  $\max_i \|x_i\|^2 \|w^*\|^2$  עדכונים.
- א"ש מרקוב:  $\mathbb{P}[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$  עבור  $t > 0, X \geq 0$ , ועבור  $X_1, \dots, X_m$  שמתפלגים i.i.d הא"ש מתקיים עבור  $\bar{X}$ .
- א"ש צ'ביצ'ב: עבור  $t > 0$  (לא בהכרח אי שלילי) נקבל  $\mathbb{P}[X - \mathbb{E}[X] \geq t] \leq \frac{\text{Var}[X]}{t^2}$ . עבור  $\bar{X}$  מקבלים א"ש דומה עם  $m$  במכנה.
- א"ש הופדינג: עבור  $a_1 \leq X_i \leq b_i$  בלתי תלויים וחסומים נקבל  $\mathbb{P}[\bar{X} - \mathbb{E}[\bar{X}] \geq \epsilon] \leq 2 \exp\left(-2m^2 \epsilon^2 / \sum (b_i - a_i)^2\right)$ . עבור משתנים i.i.d חסומים בין 0 ל 1 נקבל הסתברות של  $2 \exp(-2m \epsilon^2)$ . עבור הטלות מטבע עם הטיה  $\hat{p}$  נקבל ש  $m(\epsilon, \delta) \leq \left\lceil \frac{1}{2\epsilon^2} \cdot \log\left(\frac{2}{\delta}\right) \right\rceil$ .

פונקציה  $f : C \rightarrow \mathbb{R}$  היא קמורה אם ורק אם לכל שני וקטורים  $u, v \in C$  מתקיים  $f(\lambda u + (1 - \lambda)v) \leq \lambda f(u) + (1 - \lambda)f(v)$ . תנאי שקול לכך הוא שה  $\text{epigraph}(f) = \{(x, \beta) : f(x) \leq \beta\}$  הוא קבוצה קמורה. פונקציית נורמה היא קמורה, כל פונקציה אפיינית היא קמורה. כל צירוף לינארי של פונקציות קמורות עם סקלרים חיוביים הוא פונקציה קמורה. הרכבה של פונקציה קמורה על פונקציה אפיינית כלשהי היא פונקציה קמורה גם כן. סופרימום של מס' סופי של פונקציות קמורות היא פונקציה קמורה. פונקציה חד ממדית ודיפרנציאבילית  $f$  היא קמורה אם ורק אם הנגזרת שלה מונוטונית עולה, ואם היא דיפרנציאבילית פעמיים אז היא קמורה אם ורק אם הנגזרת השנייה אי שלילית. הפונקציה הלוגיסטית  $f(w) = \log(1 + \exp(-y \langle w, x \rangle))$  היא קמורה.

אם  $f$  היא קמורה, אזי כל מינימום לוקאלי הוא מינימום גלובאלי. בנוסף, המשיקים לפונקציה עוברים מתחת לפונקציה כלומר לכל  $u$  מתקיים  $f(u) \geq f(w) + \langle \nabla f(w), u - w \rangle$ .

$v$  נקרא סאב-גרדיאנט של  $f$  ב  $w$  אם לכל  $u$  מתקיים  $f(u) \geq f(w) + \langle v, u - w \rangle$ . פונקציה  $f$  קבוצת  $C$  של  $h$  של  $f$  בנק'  $w$  נקראת דיפרנציאל ומסומנת  $\partial f(w)$ . פונקציה  $f$  קמורה אם ורק אם לכל  $w$  יש לפחות  $SG$  אחד (כלומר  $\partial f(w) \neq \emptyset$ ). פונקציה  $f$  קמורה ודיפרנציאבילית ב  $x$  אזי  $\partial f(x) = \{\nabla f(x)\}$  אם ורק אם  $\partial f(x) = \emptyset$  הוא בנקודה  $w$  אם ורק אם  $w$  היא מינימום גלובאלי.

פונקציה  $f : C \rightarrow \mathbb{R}$  נקראת  $\rho$ -ליפשיצית אם לכל  $w_1, w_2$  מתקיים  $|f(w_1) - f(w_2)| \leq \rho \|w_1 - w_2\|$ . אם פונקציה  $f$  קמורה, אז היא  $\rho$ -ליפשיצית אם ורק אם הנורמה של כל  $SG$  היא לכל היותר  $\rho$ .

באופטימיזציה קמורה רוצים למצוא פונקציה קמורה  $f(w)$  עבור  $w \in C$ . מקרים מיוחדים: פייזיביליות -  $f$  היא פונקציה קבועה, כלומר צריך למצוא  $w \in C$ . Unconstrained minimization -  $C = \mathbb{R}^d$ .

אלג' האליפסואיד לבעיית הפייזיביליות מניח ש  $C$  חסומה מלמעלה ומלמטה וכן שניתן להפריד בין  $C$ ,  $w$  (אם  $w \notin C$ ) קיים על מישור מפרד, ובכל שלב מעדכן את האליפסואיד לפי אוב ההפרדה. מתכנס אחרי  $2d(2d+2)\log(R/r)$  צעדים לכל היותר.

מימוש של האוב הפרדה: בהינתן  $C$  נתונה בתור חיתוך של  $level - set$  של פונקציות קמורות נבדוק אם  $f_i(w) \leq 0$  לכל  $i$ . אם כן,  $w \in C$ , אחרת נחזור מינוס של  $SG$  של  $f_i(w)$ .

אם רוצים למצוא פונקציה קמורה וליפשיצית אפשר להשתמש באליפסואיד עם  $2d(2d+2)\log\left(\frac{\rho \|w^*\|}{\epsilon} + 1\right)$  צעדים.

אלג' gradient descent מתחיל מווקטור משקולות התחלתי  $w^{(1)}$  ומעדכן בכל שלב  $w^{(t+1)} = w^{(t)} - \eta \nabla f(w^{(t)})$  ( $\eta > 0$  הוא Bias-complexity tradeoff). לבסוף מחזיר  $w$  ממוצע. ב SubGD מחליפים את הגרדיאנט ב  $SG$ . SubGD מתכנס תוך  $\frac{\rho^2 \|w^*\|^2}{\epsilon^2}$  צעדים עבור פונקציה  $f$  שהיא  $\rho$ -ליפשיצית וקמורה.

על מנת להביע את הבעיה של מציאת מישור מפרד כבעיית אופטימיזציה קמורה נכתוב אותה כ  $f(w) = \max_i x_i - y_i \langle w, x_i \rangle$ . פונקציה  $f$  קמורה ו-1 ליפשיצית. SubGD מחזיר על מישור מפרד. כמעט זהה ל Batch perceptron. הגורם  $\eta$  שחסר בפרסטרון לא משנה את הסימן לכן לא חשוב.

SubGD לא תלוי ב  $d$  (רק שלב העדכון עולה לינארית עם  $d$ ) לעומת האליפסואיד שנותן דיוק הרבה יותר טוב (תלוי לוגריתמית בדיוק) אך תלוי ריבועית ב  $d^2$ . אם מדברים על מציאת על מישור מפרד אז ככל שהמדגם שלנו מופרד יותר טוב ככה נעדיף שיטות של SubGD ואם הוא מופרד גרוע נרצה לעבוד עם האליפסואיד (כי הוא תלוי רק לוגריתמית בשוליים) אך התלות היא בממד.

בעיית למידה היא קמורה אם  $\mathcal{H}$  היא קבוצה קמורה ואם לכל  $z$  מתקיים ש  $\ell(\cdot, z)$  היא קמורה.  $\mathcal{H}$  תהיה פרמטריזציה של הפונקציות. פתרון ה ERM של בעיית למידה קמורות הוא בעיית אופטימיזציה קמורה. לא כל בעיות הלמידה הקמורות הן למידות. נצטרך להוסיף אילוץ שאומר ש  $\mathcal{H}$  חסומה (ולכל  $w \in \mathcal{H}$  מתקיים  $\|w\| \leq B$ ) וש  $\ell$  ליפשיצית. ניתן להחליף את הדרישה לליפשיציות בדרישה ל"חלקות" (הנגזרת של הפונקציה היא ליפשיצית) ואי שליליות של  $\ell$ .

אם פונקציית ה  $loss$  היא לא קמורה נשתמש בפונקציית Surrogate loss שהיא כן קמורה והיא חוסמת מלמעלה את הפונקציה המקורית. למשל עבור  $1 - loss$  נגדיר  $\ell^{hinge}(w, (x, y)) = \max\{0, 1 - y \langle w, x \rangle\}$ .

באלג' SGD מעדכנים כל פעם ע"פ גרדיאנט בנקודה אחת שנבחרת אקראית מתוך  $D$  (במטרה למצוא את  $L_D$  אם רוצים למצוא את  $L_S$  בוחרים נקודה מתוך מדגם נתון) וממזערים את  $L_D(h)$  במקום  $L_S(h)$ . אין פה סתירה למשפט היסודי של הלמידה כי לא מדובר כאן על בעיית קלסיפיקציה בינארית (עם  $1 - loss$  0 שהוא לא קמור). בסוג בעיות כאלה ERM הוא לא בהכרח הפרדינמה הכי טובה. אם יש לנו בעיית למידה שהיא קמורה חסומה וליפשיצית, אם נבחר  $T \geq \frac{B^2 \rho^2}{\epsilon^2}$  צעדים של SGD עם  $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$  נקבל  $\mathbb{E}[L_D(\bar{w})] \leq \min_{w \in \mathcal{H}} L_D(w) + \epsilon$ .

קמירות לבדה לא מבטיח למידות. בעיית רגרסיה לינארית עם  $loss$  ריבועי היא לא למידה, אלא אם מניחים חסימות של  $\mathcal{X}, \mathcal{H}$ .

כשמסתכלים על סיווג מסמכי טקסט לנושא למשל, ומשתמשים בשיטת Bag of words, אם ננתח לפי חסמי VCdim נקבל Sample complexity מאוד גבוה, לעומת זאת אם ננתח לפי חסמי נורמה נקבל חסם יחסית נמוך יותר ( $R^2 \|w\|^2$ ) שלא תלוי בממד של המילון. ישנן בעיות הפוכות שבהן הממד הרבה יותר קטן מהנורמה ואז הוא עדיף.

• אלג' AdaBoost (עושה בוסטינג לדיוק) מקבל מדגם  $S$ , לומד חלש  $WL$  ומספר סיבובים  $T$ : מאתחל התפלגות אחידה  $D^{(1)}$  ובכל סיבוב  $t$  מחשב  $h_t = WL(D^{(t)}, S)$  ואז  $\epsilon_t = L_{D^{(t)}}(h_t) = \sum_{i=1}^m D_i^{(t)} \mathbb{1}_{[y_i \neq h_t(x_i)]}$  ואז מעדכן לכל  $i = [m]$  את ההתפלגות להיות  $D_i^{(t+1)} = D_i^{(t)} \exp(-w_t y_i h_t(x_i))$  (עם גורם נרמול כדי שתהיה התפלגות חוקית) כאשר  $w_t = \frac{1}{2} \log\left(\frac{1}{\epsilon_t} - 1\right)$ , ואז מחזיר את הסימן של הממוצע המשוקלל של ההיפותוזות בכל שלב. אינטואיטיבית האלג' מכריח את ההיפותוזות הבאה להתמקד בדוגמאות שהוא טעה עליהן פעם קודמת. השיגאה של  $h_t$  על ההתפלגות  $D^{(t+1)}$  היא בדיוק חצי (כך נבחרו המשוואות).

• האלג' מבטיח שאם  $WL$  הוא  $(1/2 - \gamma, \delta)$  לומד חלש, אז בסיכוי  $1 - \delta T$  מתקיים שעבור  $h_s$  שהאלג' מחזיר  $L_S(h_s) \leq \exp(-2\gamma^2 T)$ . אם ניקח  $T \geq \frac{\log(1/\epsilon)}{2\gamma^2}$  נקבל שגיאה קטנה  $\epsilon$ , ואם נבחר  $\epsilon = 1/2m$  נקבל שגיאת אימון 0 (כי היא בדידה). לרוב  $\delta$  לא מעניין כי ניתן לעשות לו בוסטינג במחיר לא גדול (בערך פ 4 יותר סיבובים של האלג').

• שגיאת האפרוקסימציה קטנה עם  $T$  ( $\approx$ סיבוביות המחלקה) ושגיאת האסטימציה גדלה עם  $T$ . ממד ה  $VC$  של מחלקת כל ההיפותוזות שהאלג' מחזיר  $L(B, T)$  (כל ההרכבות של  $T$  היפותוזות ב  $B$  שה  $WL$  מחזיר) הוא קטן מ  $O(T \cdot VCdim(B))$ . בנוסף, אם  $T = \log(m) / (2\gamma^2)$ ,  $m \geq \Omega\left(\frac{\log(1/\delta)}{\gamma^2 \epsilon}\right)$ , אז בהסתברות של לפחות  $1 - \delta$  נקבל כי  $L_{D,f}(h_s) \leq \epsilon$ .

Nonuniform learning, MDL, SRM, Decision Trees, Nearest neighbors

• סוגים שונים של ידע מוקדם יכולים להיות שהיפותוזות "קצרות" הן יותר טובות, או שדברים שגורמים דומים הם באמת דומים.

• בהינתן מחלקת היפותוזות בת מניה  $\mathcal{H}$  נגדיר  $w : \mathcal{H} \rightarrow \mathbb{R}$  כך ש  $\sum_{h \in \mathcal{H}} w(h) \leq 1$  למשל עבור  $w(h) = 2^{-|h|}$  כאשר  $|h|$  הוא האורך של המילה שמתאימה להיפותוזתה (בשפה שהיא חסרת רישות) מתקיים  $\sum_h w(h) \leq 1$ .

• אם יש לנו מחלקת היפותוזות ופונקציית משקל כ"ל  $w$ , אזי בהסתברות לפחות  $1 - \delta$  לכל היפותוזתה מתקיים  $L_D(h) \leq L_S(h) + \sqrt{\frac{-\log(w(h)) + \log(2/\delta)}{2m}}$  (בניגוד לחסם ה  $VC$  שם מחליפים את  $\log$  ב  $VCdim(\mathcal{H})$ ).

• מועדור חסם ה  $VC$  מוביל לאלג' ERM ומועדור חסם ה MDL מוביל לכלל הלמידה •  $MDL(S) \in \arg \min_{h \in \mathcal{H}} \left[ L_S(h) + \sqrt{\frac{-\log(w(h)) + \log(2/\delta)}{2m}} \right]$  כלל הלמידה הזה נותן לנו Bias-complexity tradeoff.

• MDL הוא לומד אוניברסלי, אך לומד Nonuniform. בלמידה לא יוניפורמית, מספר הדוגמאות תלוי גם ב  $h$  ולא רק ב  $\epsilon, \delta$ . למשל מחלקת כל הפונקציות החשיבות היא למידה לא יוניפורמית אך לא למידה PAC. מחלקה היא למידה לא יוניפורמית אם ורק אם היא איחוד בן מנייה של מחלקות שהן למידות PAC.

• ניתן ללמוד באופן  $NU$  את מחלקת כל הפונקציות החשיבות, אבל ב  $PAC$  לא ניתן.

•  $SRM(S) \in \arg \min_{h \in \mathcal{H}} \left[ L_S(h) + \min_{n: h \in \mathcal{H}_n} \sqrt{\frac{C^{d_n - \log(w(n)) + \log(1/\delta)}}{m}} \right]$  וכן לכל  $h \in \mathcal{H}$  מתקיים  $L_D(SRM(S)) \leq L_S(h) + \min_{n: h \in \mathcal{H}_n} \sqrt{\frac{C^{d_n - \log(w(n)) + \log(1/\delta)}}{m}}$ . כלל ה SRM הוא לומד לא יוניפורמי גנרי. מחלקת כל הפונקציות מעל דומיין אינסופי היא לא למידה לא-יוניפורמית. סיבוכיות מדגם של  $SRM : m_H \leq \min_{n: h \in \mathcal{H}_n} C^{\frac{d_n - \log(w(n)) + \log(1/\delta)}{\epsilon^2}}$ .

• מחלקת עצי ההחלטה עם  $k$  עלים היא בעלת  $VCdim = k$ . מחלקת כל עצי ההחלטה מעל  $\{0, 1\}^d$  היא מממד  $VCdim = 2^d$ . נרצה להעדיף עצים קצרים יותר ע"י שפת תיאור לעצים ואז שימוש בכלל הלמידה של MDL על מנת לחפש עץ עם  $n$  צמתים שממזער, אך בעיה זו היא NP קשה. נשתמש באלג' המקרב שעובד בצורה חמדנית שנקרא ID3 - בכל שלב הוא מפצל לפי המאפיין שממקסם את ה  $\mathbb{P}[y | -x_i] - \mathbb{P}[y | x_i]$  (Gain  $(S, i) = C(\mathbb{P}[y]) - (\mathbb{P}[x_i] C(\mathbb{P}[y | x_i]) + \mathbb{P}[-x_i] C(\mathbb{P}[y | -x_i]))$ ) וממשיך רקורסיבית עד שנמרים המאפיינים. ניתן גם לגנוס, או לבנות יער של הרבה עצים (כמעט תמיד עובד), וגם להתמודד עם Real valued features (ספים).

• בשיטת ה Nearest neighbors אין מחלקת היפותוזות. תהליך הלמידה הוא פשוט שמירת דוגמאות וביזמן הפרידקציה מחזירים את הצבעת הרוב של  $k$  השכנים הכי קרובים לדוגמא. אם מניחים  $c$ -ליפשיציות של  $\eta(x) = \mathbb{P}[y = 1 | x]$ , הפרמטר  $k$  קשור ל Bias-complexity tradeoff. בנוסף, מספר הדוגמאות גדל אקספ' עם ממד הדוגמאות. בנוסף, מספר הדוגמאות תלוי בהתפלגות  $D$ . אלגוריתם Nearest neighbors הוא Universal consistent.

• PAC חזק יותר מ Non uniform שחזק יותר מ Universal consistent.

## 6. קמירות ו SGD

• קבוצה  $C \subset V$  קמורה אם ורק אם לכל שני וקטורים  $u, v \in C$  קבוע  $\lambda \in [0, 1]$  מתקיים  $\lambda u + (1 - \lambda)v \in C$  כלומר הקטע המחבר בין הנקודות נמצא כולו בתוך הקבוצה. כל תת מרחב וקטורי  $U \subseteq V$  הוא קבוצה קמורה. חיתוך של מס' סופי של קבוצות קמורות היא קבוצה קמורה. חיבור של שתי קבוצות קמורות הוא קבוצה קמורה. הכפלה של קבוצה קמורה בסקלר כלשהו היא קבוצה קמורה. כל על מישור  $\{w, v\} = W$  הוא קבוצה קמורה.

• משפטים לא קשורים: תהי  $C$  קבוצה קמורה וסגורה ב  $V$ , ו  $u \in V$  אזי קיימת נקודה  $w \in C$  יחידה (שנקראת ההטלה של  $u$  ב  $C$ ) כך ש  $\|u - w\| \leq \|u - v\|$ . כמו כן,  $\langle u - w, v - w \rangle \leq 0$ . לכל קבוצה קמורה ונקודה שלא בקבוצה קיים על מישור מפרדי בנייה

- עבור על מישור  $L = \{v : \langle w, v \rangle = 0\}$  המרחק של נקודה מהמישור אם  $\|w\| = 1$  הוא  $d(x, L) = |\langle w, x \rangle|$ . השוליים של על מישור הוא המרחק לנקודה הכי קרובה אליו  $\min_i |\langle w, x_i \rangle + b|$  (וקטורים תומכים).
  - אלגוריתם Hard-SVM מחפש על מישור עם שוליים הכי גדולים. הבעיה הזאת היא שקולה לבעיה של למצוא את  $\|w\|^2$  כך שלכל  $i$  מתקיים  $y_i (\langle w, x_i \rangle + b) \geq 1$ .
  - בעיית Soft-SVM היא  $\min_w \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum \ell_{\text{hinge}}(y_i \langle w, x_i \rangle)$ . בעיית Soft-SVM ניתן להחליף את גורם ה loss במזעור של ממוצע  $\xi_i$  כך שלכל  $i$  מתקיים  $y_i \langle w, x_i \rangle \geq 1 - \xi_i$  כלומר  $\min_w \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum \xi_i$ .
  - באופן כללי, רגולריזציה (RLM) היא הוספת גורם שאינו תלוי במדגם למזעור  $\min_w R(w) + L_S(w)$ . Soft-SVM היא בעיית רגולריזציה. בעיה נוספת היא Ridge regression שהיא RLM עם  $\frac{\lambda}{2} \|w\|^2$ .
  - למה לעשות רגולריזציה? יכול להוריד את sample complexity - כמו ב MDL יש לנו ידע מוקדם שהיפותזות קצרות הן יותר טובות, וכן אלגוריתמים כאלה הם "יציבים". בנוסף, זה עוזר להוריד את סיבוכיות החישוב (גם לבעיות לא קמורות).
  - עבור פונקציה  $\epsilon: \mathbb{N} \rightarrow \mathbb{R}$  מונוטונית יורדת, אלגוריתם נקרא "בממוצע יציב להחלפה אחת" בקצב  $\epsilon(m)$  אם לכל הסתברות מתקיים  $\mathbb{E}[\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)] \leq \epsilon(m)$ . אם אגל' הוא "בממוצע יציב להחלפה אחת" עם קצב  $\epsilon(m)$  אז  $\mathbb{E}[L_D(A(S)) - L_S(A(S))] \leq \epsilon(m)$ .
  - משפט: רגולריזציה תיכונת  $(\lambda \|w\|^2)$  היא מייצבת - בהנחה שפונקציית ה loss קמורה וליפשיצית, אז כלל ה RLM עם רגולריזציה תיכונת היא "בממוצע יציב להחלפה אחת".  
עם  $\epsilon(m) = \frac{2\rho^2}{\lambda m}$ .
  - שיטת הקרנלים: מגדירים פונקציה  $\psi: X \rightarrow F$  כאשר  $F$  הוא מרחב מאפיינים כלשהו, ואז נלמד חצאי מרחב על  $(\psi(x_i), y_i)$ . איך נבחר את המיפוי? דורש ידע על הבעיה.
  - הקרנל של מיפוי  $\psi$  היא פונקציה שממשת מכפלה פנימית במרחב המאפיינים:  $\langle \psi(x), \psi(x') \rangle = K(x, x')$ . לפעמים קל לחשב את  $K$  בלי לחשב את  $\psi$  באופן מפורש.
  - עבור כל כלל למידה מהצורה של  $w^* = \arg \min_w (f(\langle w, \psi(x_1) \rangle, \dots, \langle w, \psi(x_m) \rangle) + \lambda \|w\|^2)$  נוכח  $w^* = \sum \alpha_i \psi(x_i)$ . לכן, אם אנחנו יודעים את  $G_{i,j} = \langle \psi(x_i), \psi(x_j) \rangle$  לחשב את  $(G\alpha)_i = \langle w, \psi(x_i) \rangle$  וכן את  $G\alpha = \alpha^T G\alpha$  וכך לקבל את בעיית האופטימיזציה:  $\arg \min_{\alpha} (f(G\alpha) + \lambda \alpha^T G\alpha)$  (שהיא בממד הרבה יותר קטן לרוב). אם אכן משתמשים בטריק הזה, כאשר רוצים לבצע פרדיקציה על  $x$  מחשבים את  $\langle w, \psi(x) \rangle = \sum_j \alpha_j K(x_j, x)$ .
  - קרנל פופולרי הוא קרנל גאוס (RBF) - לומדים פולינומים ממעלה אינסופית ע"י קרנל יעיל.
  - ע"פ תנאי מרסר - אם לכל בחירה של  $m$  דוגמאות המטריצה  $G$  היא מוגדרת אי שלילית אזי היא פונקציית קרנל חוקית. בנוסף, מכיוון של פונקציית קרנל היא מכ"פ, אזי היא חייבת לקיים את התנאים של מכ"פ (חיוביות בהחלט - לכל  $x$  מתקיים  $\langle x, x \rangle \geq 0$  ושוויון רק ב 0, לינאריות וסימטריות).
- ## 8. צברור ו Features
- במודל צברור יש קבוצת אובייקטים ופונקציית מרחק  $d$ , בעוד הפלט הרצוי היא חלוקה של קבוצת האובייקטים ל  $k$  קבוצות זהות. הפלט יכול להיות גם עץ של קבוצות כך שבכל רמה יש חלוקה גסה יותר (ברמה הכי נמוכה יש את הדוגמאות וברמה הכי גבוהה יש קבוצה אחת שמכילה את כולם).
  - בשיטות מבוססות קישורים נבנה בעצם דנדוגרם (עץ כמתואר לעיל) כאשר נתחיל מסט של  $|X|$  קבוצות וכל פעם נחבר את שתי הקבוצות שהכי דומות לפי פונקציית המרחק. יכולות להיות כמה הרחבות של  $d$  לפונ'  $d$  שמאחדת בין קבוצות ובין נקודה לקבוצה (המרחק המינימלי בין הקבוצות, המרחק הממוצע, המרחק המקסימלי וכו').
  - בשיטות מבוססות מזעור מחיר מגדירים פונקציית מחיר  $G$  אשר מגדירה עבור קלט של דוגמאות ופונקציית מרחק ביחד עם צברור מוצע את המחיר של הצברור, ואז מנסים למצוא את הפונקציה  $G$ .
  - במשפחת  $k$ -means מגדירים את  $G = \min_{\mu_i \in X'} \sum_{i=1}^k \sum_{x \in C_i} d(x, \mu_i)^2$  כאשר  $\sum_{i=1}^k \sum_{x \in C_i} d(x, \mu_i(C_i))^2 = \arg \min_{\mu \in X'} \sum_{x \in C_i} d(x, \mu)^2$  הוא המרכז של הצבר  $i$ . ניתן גם להסתכל על מרכזים שהם רק בתוך הקבוצה המקורית (או להסתכל על מרחק בלי ריבוע  $(k - \text{median})$ ).
  - בעיית  $k$ -means NP קשה, לכן משתמשים באגל' קירוב שנקרא  $k$ -means שמתחיל ממרכזים אקראיים, ובכל שלב לוקח לצבר  $i$  את כל הנקודות שהכי קרובות למרכז הזה ואז מעדכן כל מרכז להיות המרכז של כל הנקודות שמשיכות אליו.
  - איך מייצגים אובייקטים מהעולם האמיתי בתור וקטור של Features? בהנתן שיש לנו ייצוג של Features, האם זה הייצוג הכי טוב שיש?
  - בהנתן מרחב דוגמאות מממד  $d$  נרצה לייצג וקטור Features שמכיל רק  $k$  Features (עבור  $k \ll d$ ) על מנת להוריד את טעות האסטימציה, להוריד את העלות החישובית וגם את עלות השגת המידע.
  - בשיטת הפילטרים נעריך כל Feature בנפרד ע"י פונקציית מחיר כלשהי (למשל ה Feature שנותן את ה loss המינימלי הכי גבוה, או מה שמזעור את ה loss הריבועי שה שקול למקדם הקורלציה של פירסון שהוא  $\frac{|\langle v - \bar{v}, y - \bar{y} \rangle|}{\|v - \bar{v}\| \|y - \bar{y}\|}$ , או להפעיל את מקדם פירסון על הציון של ה Feature  $v$ ) וניקח את  $k$  ה Features הכי טובים.

- השיטה הכי פשוטה להתמודד עם ה tradeoff היא להתחיל מ  $m$  צעדים של חיפוש אקראי, ואז בשאר ה  $T - m$  צעדים לבחור במכונה שנתנה את ה reward הממוצע הכי טוב. בשיטה זו, אם  $m = O\left(\frac{n \log n}{\epsilon^2}\right)$  אזי לכל  $i$  מתקיים  $|\mu_i - \mu^*| \leq \epsilon$  (הופדינג וחססהאחוד).
  - בנוסף, ה regret הוא  $\text{regret} \leq 2\epsilon + \frac{n \log n}{T\epsilon^2}$ . עבור ה  $\epsilon$  הכי טוב מקבלים שה  $\text{regret} = \left(\frac{n \log n}{T}\right)^{1/3}$ .
  - טכניקה נוספת היא חיפוש  $\epsilon$ -greedy. בעצם אנו רוצים למזער את הפונקציה  $L(w) = \sum w_i = 1$  עם האילוץ  $w = [0, 1]^n$ . הבעיה היא קמורה עם אילוץים קמורים, הבעיה היא שלא יודעים לחשב את  $\nabla L$  לכן נבנה משעך בלתי מוטה לגרדיאנט. לכל וקטור התפלגויות  $p$  נבחר  $i_t \sim p$  ונקבע  $e_{i_t} = \frac{-\nabla L(w^{(t)})}{p_{i_t}}$ . נבחר את  $p = (1 - \epsilon)w^{(t)} + \epsilon \cdot \frac{1}{n}$  כלומר בסיכוי  $\epsilon$  נבחר מכונה אקראית ובסיכוי  $1 - \epsilon$  נשתמש בידע הקודם שלנו. הדבר מבטיח שהנורמה של הגרדיאנט בריבוע לא תהיה גדולה מידי וגם נקבל  $\text{regret} = \left(\frac{n}{T}\right)^{1/3}$  (רווח לגוריתמי).
  - אלג' נוסף, שדומה ל SGD הוא EXP3: מאתחלים את  $w$  להיות וקטור התפלגות אחידה ונעדכן כמו ב SGD רק באקספוננט  $- \eta \nabla L(w^{(t)})[i]$   $w_i^{(t+1)} = \frac{1}{Z_t} w_i^{(t)} \exp\left(-\eta \nabla L(w^{(t)})[i]\right)$ .
  - באלג' זה  $\text{regret} = \left(\frac{n \log n}{T}\right)^{1/2}$ . עובד גם במודל adversarial (כשטבע מתנכל לנו).
  - טכניקה נוסף היא UCB:  $\text{UCB}_i(t) = \hat{\mu}_i + \sqrt{\frac{2 \log(T)}{N_i(t)}} \geq \mu_i$ . בכל שלב נבחר את הזרוע שמקסמת את ה  $\text{UCB}_i(t)$ . ה regret בטכניקה הזאת מתנהג כמו  $\frac{\log(T)}{T}$ .
  - מודל יותר ריאליסטי הוא MDP: לכל  $t$  מתקיים  $s_{t+1} \sim \tau(s_t, a_t)$  כאשר  $\tau$  טרמיניסטי על  $S \times A$ , וכן  $r_t$  מ"מ מעל  $[0, 1]$  שתלוי רק ב  $(s_t, a_t)$  וכן  $\mathbb{E}[r_t] = \rho(s_t, a_t)$ . הנחת המרקוביות אומרת שהעבר והעבר בלתי תלויים בהנתן ההווה (ההווה מסכם את כל מה שצריך לדעת על העבר).
  - אלג' Value iteration מוצא את  $\pi$  האופטימלי כאשר  $\tau, \rho$  ידועים כבר (כאשר למשל אפשר להסיק אותם מתוך הפיזיקה). Q-Learning הוא כאשר לא יודעים את  $\tau, \rho$ .
  - באלג' Value iteration מגדירים את ה optimal value function  $V^* : S \rightarrow \mathbb{R}$  להיות מה תוחלת ה reward של הפוליסי האופטימלי בהנתן שהתחלנו ממצב  $s$ , כלומר  $V^*(s) = \mathbb{E}\left[\sum_{t=1}^{\infty} \gamma^t r_t | s_1 = s\right]$ . מכאן ניתן להסיק כי מתקיים  $V^*(S) = \max_{a \in A} [\rho(s, a) + \gamma \mathbb{E}_{s' \sim \tau(s, a)} [V^*(s')]]$ . הגורם הפנימי הוא בעצם  $Q^*(s, a) = \rho(s, a) + \gamma \mathbb{E}_{s' \sim \tau(s, a)} [V^*(s')]$  שמתאר מה יהיה ה reward אם אני נמצא כרגע במצב  $s$  ואני הולך לעשות את הפעולה  $a$  ומשם את הפעולות לפי הפוליסי האופטימלי. כלומר ה  $Q^*(s, a) = \arg \max_{a' \in A} [\rho(s, a) + \gamma \mathbb{E}_{s' \sim \tau(s, a)} [V^*(s')]]$  הוא מה שמקסם את ה  $Q^*$  בהנתן המצב הנוכחי. כלומר הפוליסי האופטימלי הוא טרמיניסטי ביחס ל  $s_t$ . האלג' אם כך פשוט מתחיל מווקטור אקראי  $V_0$  ומעדכן  $V_{t+1}(s) = \max_{a \in A} [\rho(s, a) + \gamma \mathbb{E}_{s' \sim \tau(s, a)} [V_t(s')]]$  עד להכנסות. כאשר  $V_t = V_{t+1}$  אזי הוא האופטימלי  $V^*$ .  $\|V_t - V^*\|_{\infty} \leq \gamma \|V_0 - V^*\|_{\infty}$ . כלומר הוא מתכנס מהר יחסית. ההוכחה: מגדירים אופרטור  $T^* : \mathbb{R}^{|S|} \rightarrow \mathbb{R}^{|S|}$  להיות  $V_{t+1} = T^*(V_t)$  ואז מראים שהאופרטור הוא אופרטור כיווץ (כלומר לכל שני וקטורים  $u, v$  מתקיים  $\|u - v\|_{\infty} \leq \gamma \|T^*(u) - T^*(v)\|_{\infty}$ ) ואז משתמשים במשפט נקודת השבת בנג. כעת לאחר שהתכנסנו ל  $V^*$  נוכל לחשב את  $Q^*$  ומשם את  $\pi^*$ .
  - שיטה נאיבית בהנתן שלא יודעים את  $\tau, \rho$  היא לשערך את  $\tau, \rho$  בעזרת "חיפוש אקראי" בעולם, כלומר פוליסי אקראי.
  - באלג' Q-Learning מנסים לשערך את פונקציית  $Q^*$  בלבד (model-free) ע"י ניסיון לפתור את משוואת בלמן המתאימה לה:  $Q^*(s, a) = \rho(s, a) + \gamma \mathbb{E}_{s' \sim \tau(s, a)} \max_{a'} Q^*(s', a')$ . בהנתן  $(s_t, a_t, s_{t+1}, r_t)$  נגדיר  $\delta_{s_t, a_t}(Q) = Q_t(s_t, a_t) - (r_t + \gamma \max_{a'} Q_t(s_{t+1}, a'))$  ונעדכן  $Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) - \eta_t \delta_{s_t, a_t}(Q_t) \mathbf{1}_{[s=s_t, a=a_t]}$ . בניסיון להתכנס למשוואת בלמן. עובד פחות טוב מ Value iteration. בחירת  $a_t$  בכל שלב יכולה להעשות באותן טכניקות כמו Multi armed bandit.
  - האלג' Q-Learning בצורה יותר מפורשת: מאתחלים  $Q_0 = 0$  ואז לכל רביעיה  $(s_t, a_t, r_t, s_{t+1})$  מעדכנים  $Q_{t+1}(s_t, a_t) = (1 - \eta_t) Q_t(s_t, a_t) + \eta_t (r_t + \gamma \max_a Q_t(s_{t+1}, a))$ .
  - הבעיה בשני האלג' האלה היא שהסיבוכיות עולה אקספ' עם המימד (גם sample complexity וגם סיבוכיות החישוב). נוכל לעשות אפרוקסימציה לפונקציות 12. כללי  $Q$  ה "ע" מחלקת היפותזות שתאופיין בצורה פרמטרית עם פרמטר  $\theta$  ואז נוכל לבצע את האלג' בצורה יעילה חישובית (למשל עם GD). Deep-Q-Learning. למשל זה בעצם לכת את מחלקת ההיפותזות להיות רשת נוירונים. למשל DeepMind לקחו את הפונקציה  $\theta_\theta : S \rightarrow \mathbb{R}^{|A|}$  להיות רשת שמקבלת כקלט ייצוג  $d$ -ממדי (ממש) של המצב הנוכחי ועם שכבת פלט בגודל  $|A|$  מוציאה את הערך של  $Q$  לכל אחת מהפעולות האפשריות.
  - פתרון אחר לבעיית הממד הוא באמצעות temporal abstraction המשתמש במנגנון שנקרא אופציות: כל פעם בוחרים אופציה (למצוא את הטלפון) שבתורה קוראת לאופציה נחתה בתהליך היררכי, עד שהיא מסתיימת ואז חוזרים לתהליך העליון, ועוברים לאופציה הבאה. כלומר מחלקים את הבעיה לתתי בעיות, מייצרים מבנה היררכי על מרחב המצבים ומתמודדים כל פעם רק עם חלק קטן יותר מהמרחב.
- 11. הורדת ממד**
- לוקחים מידע שהוא בממד גבוה ומורידים לממד נמוך: יכול להקטין את זמן האימון וזמן הריצה, מוריד את טעות האסטימציה. טכניקות לינאריות - מכפילים  $x \in \mathbb{R}^d$  במטריצה  $W \in \mathbb{R}^{n \times d}$ ,  $n < d$ .
  - בשיטת PCA משתמשים לרוב לפני שלומדים על המידע, כשלב עיבוד מקדים. בשיטות הורדת ממד נרצה שלאחר שחזרו של  $\tilde{x} \approx x$  נקבל  $\tilde{x} \approx x$  PCA מניח שגם השחזור הוא לינארי. בעצם רוצים למצוא מטריצות  $W, U$  המתאימות ל  $A = \sum_{i=1}^m x_i x_i^T$  הפתרון של  $\arg \min \sum_{i=1}^m \|x_i - UWx_i\|^2$  (מטריצת השחזור) היא מטריצה שהעמודות שלה הן הווקטורים  $u_1, \dots, u_n$  שהם הווקטורים העצמיים המתאימים לע"ע הגדולים ביותר של  $A$ .
  - בהוכחה מראים קודם כל ש  $W = U^T$  והעמודות של  $U$  הן בסיס א"ל ל  $S = \text{Im}(WU)$  ע"י כך שרואים כי הנקודה  $\tilde{x}$  שהכי קרובה ל  $x$  על  $S$  היא  $VV^T x$  כאשר העמודות של  $V$  הן בסיס א"ל, לכן אם ההנחה לא מתקיימת נוכל להחליף את  $W, U$  ורק לשפר את מצבנו. אנליזה של האלג' מראה לנו גם ששיגאת השחזור שלנו היא בדיוק סכום הע"ע שלא לקחנו.
  - ברוב המקרים מורידים את הממוצע של הדוגמאות לפני שמפעילים עליהן PCA.
  - במקרה שבו  $d$  גדול ממש עושים את הטרנספורמציה הבאה: מגדירים  $X$  להיות מטריצה שהשורות שלה הן הדוגמאות, ואז  $A = X^T X$ . נגדיר את  $B = XX^T$  שהיא מממד  $m \times m$  שהוא יותר קטן מ  $d \times d$ . כעת אם  $Bu = \lambda u$  אז נקבל כי  $A(X^T u) = \lambda(X^T u)$  לכן  $\lambda(X^T u) = \frac{X^T u}{\|X^T u\|}$  הוא ו"ע של  $A$  עם אותו ע"ע. בגלל ש  $B_{ij} = \langle x_i, x_j \rangle$  תלוייה רק במכ"פ ניתן להחליף אותה בפונקציית קרנל. במקרה שמשתמשים ב  $B$  במקום  $A$  בזמן הריצה הוא  $O(m^3 + m^2 d)$ .
  - לפעמים הורדת ממד עוזרת גם ליצירת צברים (כמו בדוגמא עם התמונות של האנשים).
  - ב Random projections לא אכפת לנו משחזור מוצלח, אלא מעניין אותנו לשמור מרחקים, כלומר שהיחס בין המרחק לפני למרחק אחרי יהיה בערך 1  $\frac{\|Wx_i - Wx_j\|}{\|x_i - x_j\|} \approx 1$ . נייער מטריצה  $W_{ij} \sim N(0, 1/n)$  ונראה מה יקרה להורדת ממד שתבצע איתה. התוחלת  $\mathbb{E}[\|Wx\|^2] = \sum_{i=1}^n \mathbb{E}[\langle (w_i, x) \rangle^2] = \sum_{i=1}^n x^T \mathbb{E}[w_i w_i^T] x = \|x\|^2 \left(\frac{1}{n} I\right) x = \|x\|^2$  יש משפט שאומר שאם  $n \geq \frac{\log(|Q|) - \log(\delta)}{\log(\epsilon^2)}$  (סדר גודל, עד קבועים) אז הסיכוי שהורדת ממד מקרית תעוות את המרחקים ביותר מ  $\epsilon$  קטן מ  $\delta$ . באופן כללי אם יש לנו  $n$  דוגמאות בממד מגודל  $d$  אז בפועל הדוגמאות יושבות בתת מרחב מממד  $n$ .
  - הלמה של ג'ונסון ולינדרשטראוס: תהי  $Q$  קבוצה סופית ב  $\mathbb{R}^d$ ,  $\delta \in (0, 1)$  ו  $n$  מס' טבעי כך ש  $\epsilon = \sqrt{\frac{6 \log(2|Q|/\delta)}{n}} \leq 3$  אז בסיכוי לפחות  $1 - \delta$  על בחירת מטריצה מקרית  $W \in \mathbb{R}^{n, d}$  עם  $W_{i,j} \sim N(0, 1/n)$  מתקיים:  $\max_{x \in Q} \left| \frac{\|Wx\|^2}{\|x\|^2} - 1 \right| < \epsilon$ . כלומר הטלה מקרית משמרת מרחקים עבור קבוצה סופית של וקטורים. בנוסף, אם הוקטורים יושבים בתוך כדור יחידה, גם המכפולות הפנימיות שלהם נשמרות.
  - בשיטת Compressed sensing מניחים כי  $U \approx U \alpha$  עבור  $U$  אורתונורמלית כך שמספר האינדקסים השונים מ0 ש  $U \alpha$  הוא  $d \gg d$ . נשמור את המידע אם כך ע"י מציאת  $\alpha = U^T x$  ושמירת כל האלמנטים שהם לא 0 בווקטור  $\alpha$  ואת האינדקסים שלהם (דורש בערך  $s \log(d)$  ביטים).
  - Compressed sensing מבטיחים שחזור מושלם לכל הוקטורים שהם  $O\left(\frac{n}{\log(d)}\right)$  דלילים. PCA לעומת זאת מבטיח שחזור מושלם אם הדוגמות יושבות בדיוק בתמ"ו מממד  $n$  (מוצא תמ"ו מממד  $n$  שהכי קרוב לדוגמאות). אם עושים בדיוק  $n$  compressed sensing יכול להיות שהם יכשלו הדוגמה הזאת, אך אם נעשה  $n \log(d)$  compressed sensing - מכיוון שיש בסיס שבו המידע הזה הוא  $n$  דליל, לכן אם ניקח קצת יותר, כלומר  $n \log(d)$  אז נצליח. באופן פרקטי PCA דליל מוצלח מ compressed sensing ברוב המקרים.
  - Compressed sensing יותר טוב בדוגמה לעיל מכיוון שבדוגמה זו השחזור הוא לא לינארי, כי זאת אחת ההנחות של PCA (בשחזור לינארי PCA הוא אכן האופטימלי).
  - מטריצה  $W$  היא RIP  $(\epsilon, s)$  אם לכל  $x \neq 0$  כך ש  $\|x\|_0 \leq s$  מתקיים  $\left| \frac{\|Wx\|_2^2}{\|x\|_2^2} - 1 \right| \leq \epsilon$ . משפט: עבור  $\epsilon < 1$  תהי  $W$  מטריצה שהיא RIP  $(\epsilon, 2s)$ , וקטור  $x$  ו  $\|x\|_0 \leq s$ ,  $y = Wx$  ו  $\|y\|_0 = v$  ו  $y = v$  ו  $\|y\|_0 = v$  אזי  $\tilde{x} = x$ . כלומר לפי המשפט נסיק כי תכונת ה RIP משמרת את השחזור האופטימלי לכל הווקטורים הדלילים. בנוסף, אם נוסיף דרישה  $\epsilon < \frac{1}{1+\sqrt{2}}$  אז גם נורמטת 1 מקיימת את התכונה ל  $\tilde{x}$  (כלומר אם  $\tilde{x}$  הוא הווקטור עם הנורמטת 1 הכי קטנה מבין אלה שמשוחזרים ל  $x$ ). את הבעיה עם נורמטת 1 ניתן לכתוב כבעיית LP. כלומר אפשר לשחזר את כל הווקטורים הדלילים בצורה יעילה אם מורידים אותם מקרית ל  $s \log d$  מדידות.
  - חוקי log:  $\log(xy) = \log(x) + \log(y)$ ,  $\log_b(x) = \frac{\log_c(x)}{\log_c(b)}$ .
  - מטריצה סימטרית ממטית  $A$  נקראת PSD אם לכל  $v \in \mathbb{R}^d$  מתקיים  $v^T A v \geq 0$ , או אם קיימת  $B$  כך ש  $A = BB^T$ , או אם כל הע"ע א"ש.
  - bayes optimal classifier בהנתן  $D$  הוא  $\mathbb{P}[y = 1|x] \geq 0.5$   $f_D(x) = \begin{cases} 1 & \text{if } \mathbb{P}[y = 1|x] \geq 0.5 \\ 0 & \text{o.w.} \end{cases}$
  - SVM עם קרנל RBF יכול להביע את כל הפונקציות מ  $\{\pm 1\}^d \rightarrow \{\pm 1\}$  מבמבחן: עבור  $r$  מחלקות היפותזות, ו  $d \geq 3$  הוא ה VCdim המקסימלי, אזי  $\text{VCdim}(\bigcup_{i=1}^r \mathcal{H}_i) \leq 4d \log(2d) + 2 \log(r)$  אם  $r = 2$  נקבל  $2d + 1$ .
  - $\text{PAC} \leq \text{NUL} \leq \text{consistent}$