

Principal Components Analysis (PCA)

נוריד את הדאטה ל- $k$  מימדים כך שנישאר עם  $Var$  מקסימלי. הוקטורים העצמיים מאונכים זה לזה. נקודה  $x_j \in \mathbb{R}^d$  בהטלה ל- $k$  מימדים תיוצג ע"י:

$$error = \left\| x_j - \sum_{i=1}^k \alpha_{ji} e_i \right\|$$

המטרה: מזער את  $\left\| x_j - \sum_{i=1}^k \alpha_{ji} e_i \right\|$ .

1. חשב תוחלת  $\hat{\mu}, \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$  נגדיר  $z_i = x_i - \mu$

2. חשב את  $S = \sum_{j=1}^n z_j z_j^T$  ומצא  $k$  וקטורים עצמיים עם הערכים העצמיים הגדולים ביותר

3. נוריד מימד ע"י  $Y = E^T Z$  ( $E = [e_1, \dots, e_k]$ )

האלגוריתם הוא unsupervised, לא מתייחס ל-labels. שחזר ע"י כפל ב- $E^T$  (והוספת  $\mu$ )

Logistic Regression

$y(x) = \text{sign}(w^T x + w_0)$  נרצה לגרום לפונקציה להיות יותר חלקה וגזירה, לכן נשתמש "sigmoid":  $\sigma(z) = \frac{1}{1+e^{-z}}$ . יש חשיבות למרחק של הנקודה מוקו ההפרדה. למציאת ה- $w$  האופטימלי נמוזער את פונקציית ההפסד עם gradient descent, לכל  $j \in [1, d]$ :

$$w_j^{t+1} = w_j^t - \eta \sum_i x_i^j (\mathbb{P}(c = 1 | x_i; w) - y_i)$$
$$\ell(w) = \sum_{i=1}^n \ln \left( 1 + e^{-y_i a^T x_i} \right) = \frac{1}{n} \sum_{i=1}^n y_i \log h(x_i) + (1 - y_i) \log (1 - h(x_i))$$

Linear discriminant analysis (LDA)

נמצא הטלה לקו שתשמר את ההפרדה (נעדין יש צורך לסווג לאחר מכן) נרצה למקסם את מרחק תוחלות ההטלה ולמזער את השוניות, נמקסם את  $J(v) = \frac{(\bar{\mu}_1 - \bar{\mu}_2)^2}{S_1^2 + S_2^2}$ . לא מובטחת הפרדה מוחלטת, אבל נקבל הפרדה טובה ביותר.  $v$  שממקסם את  $J$  מקיים  $S_B v = \lambda S_W v$  כאשר  $S_B = S_1 + S_2$ .

$$S_j = \sum_{x_i \in c_j} (x_i - \mu_j) (x_i - \mu_j)^T$$
$$S_B = (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T$$

אם  $S_W$  הפיכה נמצא ישירות:  $v = S_W^{-1} (\mu_1 - \mu_2)$ . לעומת PCA, משנה רק הכיוון של הוקטור שמצאנו ולא מיקומו במרחב. LDA ממירה מכל מספר פיצ'רים למימד אחד. האלגוריתם הוא supervised. MDA - הכללה עבור  $c \geq 2$  מחלקות, ניתן להוריד עד  $c - 1$  מימדים, נסווג למחלקה עם ערך הכי גבוה.

$\epsilon$ -Representitives

קבוצת אימון  $S$  היא  $\epsilon$ -ייצוגית אם לכל היפותזה  $h$ :  $|L_S(h) - L_D(h)| \leq \epsilon$ . כלומר,  $S$  מייצגת את המציאות עד כדי  $\epsilon$ .  $\mathcal{H}$  מתכנסת באופן אחיד אם בהסתברות לפחות  $1 - \delta$  מדגם  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  יהיה  $\epsilon$ -מייצג לכל היפותזה  $h$ . אם  $\mathcal{H}$  מקיימת תכונה זו אז היא למידה ב- $\text{Agnostic PAC}$  עם סיבוכיות דגימות  $m_{\mathcal{H}}^{\text{ERM}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{\text{UC}}(\frac{\epsilon}{2}, \delta)$  וגם  $\text{ERM}_{\mathcal{H}}$  יעבוד טוב עבורה.

$k$ -nearest neighbors (KNN)

לא פרמטרי, ודיסקרימינטיבי. ב-KNN נשערך את ה-posterior. נקבע ערך  $k$  ונמצא  $V$  שיכיל בתוכו  $k$  נקודות. יכול להתכנס לאותה תוצאה כמו KNN.PW לא תחשב פונקציית צפיפות, וקו ההחלטה לא גזיר.

$$p(c_i | x) = \frac{k_i}{k}; \quad p(x, c_i) \approx \frac{k_i/n}{\sqrt{V}}$$

עבור 0-1 loss, כלומר MAP, המחלקה עם ה- $k_i$  הגדול ביותר תיבחר. סיבוכיות גדולה מאוד,  $O(ndk)$  ומספר הפרמטרים גדל יחד עם הקלט.

- $k$  קטן מדי יגרום ל- $\text{decision boundaries}$  קופצניים ולא מאוזנים (Var גבוהה)
- $k$  גדול מדי יגרום לגבולות החלטה חלקים מדי (bias גבוה)

להסתכל על היחס בין הצירים בגרף! בנוסף, אין תהליך אימון.

Naive Bayes

מניחים מאפיינים ב"ת בהיתן  $y$ .  $\mathbb{P}(c_i | X) = \frac{\mathbb{P}(c_i) \times \mathbb{P}(X | c_i)}{\mathbb{P}(X)}$  אם באמת ב"ת, אופטימלי:

$$\mathbb{P}(w_j) = \frac{|w_j|}{|D|}; \quad c^* = \arg \max_{c_i} \mathbb{P}(c_i) \times \prod_j \mathbb{P}(X_j | c_i)$$

laplace smoothing - כאשר במחלקה כלשהי לא מופיעה דגימה נחליק את הנתונים:

$$\mathbb{P}(X = x | w_k) = \frac{\#x \in w_k + 1}{|w_k| + d}; \quad d = \# \text{פיצ'רים}$$

Bayesian Estimation:  $\theta$  הוא RV עם התפלגות משלו.

$$\mathbb{P}(X | D) = \int \mathbb{P}(X | \theta) \times \mathbb{P}(\theta | D) d\theta; \quad \mathbb{P}(\theta | D) = \frac{\mathbb{P}(D|\theta) \times \mathbb{P}(\theta)}{\int \mathbb{P}(D|\theta) \times \mathbb{P}(\theta) d\theta}$$

כאשר  $n \rightarrow \infty$  Bayesian Estimation ישאר ל-MLE. לא תמיד בייס יותר טוב מ-MLE, תלוי ב-prior ובמספר הדגימות.

Parzen Windows

לא פרמטרי, וגנרטיבי. ב-Parzen Windows נשערך את ה-likelihood. נסמן

$$p(x) = p(x | c_i) \approx \frac{k/n}{V}$$

ב-PW קובעים את החלון, קוביה  $d$ -מימדית עם צלע  $h$ :  $V = h^d$ . נבדוק האם נקודה  $x$  נמצאת בחלון (region) נשתמש בפונקציית האינדיקטור:

$$\Phi\left(\frac{x_i - x}{h}\right) = \begin{cases} 1, & \frac{|x_i - x|}{h} \leq \frac{1}{2}, k \in \{1, \dots, d\} \\ 0, & \text{else} \end{cases}$$

$$p(x) \approx \frac{k/n_j}{V} = \frac{1}{n} \sum_{i=1}^{n_j} \frac{1}{h^d} \Phi\left(\frac{x_i - x}{h}\right); \quad k = \sum_{i=1}^n \Phi\left(\frac{x - x_i}{h}\right)$$

נדרוש  $\Phi \geq 0$  וגם  $\int_{\mathbb{R}} \Phi(u) du = 1$ . נרצה לבחור  $h$  שיגרום לצפיפות להיות חלקה.

- $h$  קטן מדי יגרום לפונקציה להיות מאוד קופצנית, כל דגימה תשפיע קצת מדי
- $h$  גדול מדי יגרום לפונקציה להיות חלקה מדי, כל דגימה תשפיע יותר מדי

נסוג את המחלקה שממקסמת את  $p(x)$ . קשה להעריך את ה-אופטימלי. נוכל לחלק את הצפיפות ע"י  $\varphi$  גאוסיאנית, וכך הצפיפות הסופית תהיה ממוצע גאוסיאנים (שמתפלגים  $\mathcal{N}(0, 1)$ ).

מסווג MLE: החיזוי הוא  $\arg \max_{c_i} \mathbb{P}(x | c_i)$

מסווג MAP: החיזוי הוא  $\arg \max_{c_i} \mathbb{P}(x | c_i) \cdot \mathbb{P}(c_i)$

ה-prior הוא  $\mathbb{P}(c_j)$ , ה-likelihood היא  $\mathbb{P}(x | c_j)$

ה-posterior הוא  $\mathbb{P}(c_j | x)$  (מזה אכפת לנו)

סיכון מותנה:  $R(a_i) = \sum_{j=1}^m \mathbb{P}(c_j | x) \times \lambda(a_i | c_j)$  סיכון כולל:  $R(a_i) = \int \mathbb{P}(a(x) | x) \times \mathbb{P}(x) dx$  bayesian decision rule - נרצה למזער את הסיכון המותנה עבור 0-1 loss. כלל החלטה ביאסיאני שקול ל-MAP (אם אין priors (MLE

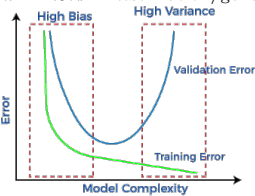
$$\mathbb{P}(err | \ell) = 1 - \mathbb{P}(pred | \ell); \quad \mathbb{P}(err) = \int_{-\infty}^{\infty} \mathbb{P}(err | x) \times \mathbb{P}(x)$$

ב-MLE נשערך את הפרמטר  $\theta$  שיקרב באופן הטוב ביותר את הדגימות להתפלגות. מניחים כי המאפיינים ב"ת. גזירים ומשווים  $\nabla = 0$ .  $\hat{\theta} = \arg \max \ln \mathbb{P}(x_1, \dots, x_n | \theta)$  משערך  $\hat{\theta}$  הוא unbiased אם  $\mathbb{E}(\hat{\theta}) = \theta$ . בנוסף, נשאף לשונות נמוכה.  $Uni(0, \theta) : \hat{\theta} = \max_i X_i$

$$\mathcal{N}(\mu, \sigma) : \hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i, \hat{\sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^2$$

Error Decomposition

השגיאה  $\min_{h \in \mathcal{H}} L_D(h)$  תלויה רק ב- $\mathcal{H}$ . כאשר היא קטנה ומניחים עליה יותר - bias גבוה - underfitting. אם מגדילים את סיבוכיות  $\mathcal{H}$  שגיאת האימון תקטן, החל ממקום מסוים השגיאה בפועל תעלה - variance גבוה - overfitting. שגיאות: ממודל לא מאומן למאומן - optimization. ממודל מאומן למודל הטוב ביותר במחלקה - estimation/generalization. מהטוב ביותר למחלקה לאופטימום - modeling.



Perceptron

סיווג לא נכון כאשר  $a^t y_i < 0$ . כל עוד קיימת דוגמא שמסווגת לא נכון, עבור על כל דוגמא ועדכן במידה ומסווגת לא נכון:  $a^{k+1} = a^k - \eta^k X_i$ . אם הדאטה לא ניתן להפרדה לינארית נרצה  $\eta^k \rightarrow 0$  כאשר  $k \rightarrow \infty$ . למשל,  $\eta^k = \frac{1}{k}$ . אין לנו הבטחה שהתהליך יעצור במקום טוב במידה והדאטה לא ניתן להפרדה לינארית. batch rule: בכל פעם נעדכן עם כל הדוגמאות שמסווגות שגויה,  $a^{k+1} = a^k - \eta^k \nabla J(a_k)$ . כך התהליך חלק יותר, לעומת ה-single sample rule/stochastic בו יש הרבה קפיצות. הפסד:  $J(a) = \sum_{i=1}^n \max\{0, -a^t y_i X_i\}$ .

Mean Squared Error (MSE)

מאתחלים את הוקטור  $b$  (לרוב באחדות), נוסף מימד אחד לכל דוגמא עם הערך 1, ונהפוך את הערכים לכל דוגמא במחלקה השנייה. אם  $Y$  הפיכה נחשב את  $a = bY^{-1}$ . ננסה לחשב את  $a = (Y^t Y)^{-1} Y^t b$  אם נכשלו נעדכן בכל פעם:

$$J(a) = \sum_{i=1}^n (a^t y_i - b_i)^2; \quad a^{k+1} = a^k - \eta^k Y^t (Y a^k - b)$$

אם  $\eta^1 = \eta^k$  אז  $a^k$  יתכנס ל- $a$  כך  $b^t Y a = Y^t Y a$ . לעומת פרסטרוון, תמיד מתכנס אך רגיש ל-outliers.

## Clustering

למידה unsupervised, למזער מרחק בתוך cluster ולמקסם בין cluster-ים. פונקציית דמיון  $s(x_i, x_k)$  גדלה ככל ש- $x_i$  ו- $x_k$  דומים. מרחק  $d$  להפך.  $d$  פופולריות: נורמה-2 (אוקלידי), נורמה-1 (מנהטן) או נורמה- $\infty$  (צ'בישב). פונקציית דמיון cosine:  $s(x_i, x_j) = \frac{x_i^T x_j}{\|x_i\| \cdot \|x_j\|}$ . אדישה ל-scaling (אכפת רק מהאווית). דרישות על פונקציית מרחק: סימטריות,  $d(a, a) = 0$  ואי-שוויון המשולש.

לא טוב לנרמל את הדאטה. נעדיף מספר דגימות אה בכל cluster. כשהקבוצה לא סגורה וחסומה נתקשה להפריד.  $J_{SSE} = \sum_{i=1}^k \sum_{x \in D_i} \|x - \mu_i\|^2$ . SSE עובד כאשר הדגימות ניתנות לעטיפה ע"י אליפסות כאלה, כשה-cluster-ים ניתנים להפרדה. ב-k-means ננסה למצוא מספר cluster-ים שימזערו את  $J_{SSE}$ . נבחר  $k$  ו- $k$  מרכזים (אפשר להגריל). כל עוד החלוקה משתנה:

- עבור כל כל דגימה ושייך אותה ל-cluster עם המרכז הכי קרוב
- עבור כל cluster נחשב את הממוצע והפוך אותו למרכז החדש

השגיאה קטנה בכל שלב, נשים לב שניתן להיתקע או לתת תוצאות שגויות כתוצאה מאתחול לא אופטימלי של המרכזים.

## Decision Trees

לא פרמטרי, צמצום האנטרופיה אם נבחר בפיצ'ר A:

$$Gain(A) = i(N) - \sum_{v \in Val(A)} \frac{|S_v|}{|S|} \cdot i(N_v)$$

מוניעת overfitting ע"י עצירה לפני העץ המלא ואפשר שגיאות אימון. בנוסף, ניתן לגדל עד הסוף ואז להעיק חלק מהצמתים (pruning). העץ הטוב ביותר נבחר ע"י tradeoff בין גודל העץ ושגיאת האימון. ב-pruning נוריד צמתים שתורמים מעט, באופן חומך החל מאבות העלים. עבור כל צומת:

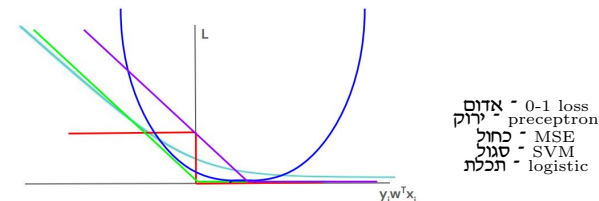
- הורד את הצומת ותת העץ שלה, החלף בעלה עם class נפוץ ביותר
- נלך ל-test ובדוק את השגיאה, נשאר רק אם הוא טוב מאחרים

נסיים כאשר הגענו למצב בו כמעט ואין עוד שיפורים, ונבחר את העץ עם השגיאה הקטנה ביותר (כמה עצים עם אותה שגיאה - ניקח את הקטן ביותר).

## Ensemble Methods

הרבה מסווגים פשוטים, כל אחד עם שגיאה קטנה מ- $1/2$  על אותה קבוצת אימון. נוכל ליצור קבוצות זרות, להגריל דוגמאות לכל מסווג או להשתמש בכל הדאטה אבל עם אתחול משקלים שונה.

ב-Bagging, בהינתן דאטה לאימון  $S$ , נדגום דאטה לכל מסווג ונאמן, בהסקה ניקח את המחלקה הנפוצה בקרב המסווגים (ברגרסיה אפשר ממוצע). טוב כאשר המסווג רגיש לרעש, משפר ביצועים וחלק יותר. ב-Random Forests נעשה Bagging על הדאטה וגם על הפיצ'רים. בבניית כל עץ נגריל את המאפיינים שיוכל להשתמש בהם ונבנה את העץ באמצעותם. עדיף מ-Bagging על עצי החלטה - מפחית תלות בין עצים שונים.



## Probably Approximately Correct (PAC)

נניח שקיים מסווג מושלם  $f$ , נגדיר את השגיאה בתור:

$$L_{D,f}(h) = \mathbb{P}_{x \sim D}[h(x) \neq f(x)]$$

נניח כי הדגימות i.i.d וכי  $y_i = f(x_i)$  לכל  $x_i$ .  $L_{D,f}(f) = 0$  וכך  $L_{D,f}(h) \leq \epsilon$  או אפילו  $L_{D,f}(h) = 0$ : No Free Lunch. יהיו  $0 < \delta < 1$ ,  $0 < \epsilon < 1/2$ , אזי לכל אלג' למידה A וקבוצת אימון בגודל  $m$  קיים עולם  $D, f$  כך ש- $\delta \geq \mathbb{P}(L_{D,f}(A(S)) \geq \epsilon)$ . אם אין מידע נוסף, לא נוכל לעשות משהו טוב יותר מאשר ניחוש. לא ידוע שום דבר על  $\mathcal{H}$  - ה-bias הוא 0 והשוונות ענקיות. עבור מחלקת היפותאות סופית  $\mathcal{H}$ ,  $ERM_{\mathcal{H}}(S)$  ימזער את:

$$L_S(h) = \frac{1}{m} |\{i \mid h(x_i) \neq y_i\}|$$

כפלט נקבל היפותאה כלשהי, אם  $f \in \mathcal{H}$  אז נקבל את  $f$ :  $L_S(f) = 0$ .  $L_S(h)$  היא שגיאת אימון ו- $L_{D,f}(h)$  היא השגיאה בפועל. אם  $\mathcal{H}$  מחלקת היפותאות לינארית אז  $\mathbb{E}_{S|x}[L_S(h)] = L_{D,f}(h)$ .

מחלקת היפותאות  $\mathcal{H}$  היא למידה ב-PAC אם קיימת פונקציה  $m_{\mathcal{H}} : (0, 1) \rightarrow \mathbb{N}$  ואלג' למידה כך שלכל  $D, f : \mathcal{X} \rightarrow \{0, 1\}$ ,  $0 < \delta < 1$ , אם נאמן את האלג' על  $m_{\mathcal{H}}(\epsilon, \delta) \geq m$  דוגמאות מאותה התפלגות ו- $i.i.d$  נקבל היפותאה  $h$  כך ש- $\mathbb{P}(L_{D,f}(h) \leq \epsilon) \geq 1 - \delta$ .  $m_{\mathcal{H}}$  היא סיבוכיות הדגימות של  $\mathcal{H}$ . אם  $\mathcal{H}$  סופית, אז  $\mathcal{H}$  למיד ב-PAC עם סיבוכיות דגימות של  $\frac{\log(\mathcal{H}/\delta)}{\epsilon}$  או  $\frac{\log(\mathcal{H}/\delta)}{\epsilon}$  אם נשתמש ב- $ERM_{\mathcal{H}}$ .

ב-Agnostic PAC לא מניחים שקיים מסווג מושלם או שהוא ב- $\mathcal{H}$ . מחלקה למידה ב-Agnostic PAC ביחס לקבוצה  $Z = \mathcal{X} \times \mathcal{Y}$  ופונקציית הפסד  $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}^+$  אם קיימת פונקציה  $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$  ואלג' A כך שלכל  $0 < \epsilon, \delta < 1$ ,  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$  ו-

$$D^m(\{S \in Z^m \mid L_D(A(S)) \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon\}) \geq 1 - \delta$$

$$L_D(h) = \mathbb{P}_{(x,y) \sim D}(h(x) \neq y) = \mathbb{E}_{z \sim D}[L(h, z)]$$

אם  $\mathcal{H}$  סופית, פונקציית הפסד בתחום  $[0, 1]$  אז היא למידה ב-Agnostic PAC עם  $ERM_{\mathcal{H}}$ , עם סיבוכיות דגימה של  $\left\lceil \frac{2 \log(2|\mathcal{H}|/\delta)}{\epsilon^2} \right\rceil$  או  $m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{2 \log(2|\mathcal{H}|/\delta)}{\epsilon^2} \right\rceil$ . עבור  $\mathcal{H}$  אינסופית  $m_{\mathcal{H}}(\epsilon, \delta) \in O\left(\frac{1}{\epsilon^2} (VCdim(\mathcal{H}) + \ln(\frac{1}{\delta}))\right)$ .

## Linear Regression

בניגוד ל-classification, נחזה ערך רציף  $y = w^T x + b$ .

$$J = \sum_{i=1}^n \left( y_i - (w^T x_i + b) \right)^2$$

נגזור ונקבל:

$$\sum_{i=1}^n x_i (y_i - w^T x_i - b) : \sum_{i=1}^n (y_i - w^T x_i - b)$$

אם המשתנים תלויים, מימד גבוה מדי או מעט מדי דוגמאות אימון נקבל משקלים עצומים. נפתור עם רגולריזציה: ridge ו-lasso. ב-ridge נגדיר:

$$J(W) = \sum_{i=1}^n \left( y_i - b - w^T x \right)^2 + \lambda \|w\|_2^2, \lambda \geq 0$$

$\lambda$  גדול ימנע מהמשקלים להתפוצץ. בנוסף,  $w_{ridge} = \left( X^T X + \lambda I \right)^{-1} X^T y$ . הפיכה.  $X^T X$  ה-ridge הפחות חשובים יקטנו אך עדיין יהיו בעלי משקל כלשהו, ב-lasso הם יתאפסו לגמרי. ל-ridge נוסחא סגורה ול-lasso אין.  $J_{lasso}(w) = \sum_{i=1}^n \left( y_i - b - w^T x \right)^2 + \lambda \|w\|_1$ .

## VC-dimension

עבור מחלקת היפותאות  $\mathcal{H}$ , ה- $VCdim(\mathcal{H})$  שלה הוא גודל הקבוצה הגדולה ביותר של דגימות שיכולה לקבל את כל התוויות האפשריות ע"י מסווגים ב- $\mathcal{H}$ . אם נרצה להוכיח ש- $VCdim(\mathcal{H}) = d$  נראה כי:

- קיימת קבוצה בגודל  $d$  ש- $\mathcal{H}$  מנתצת
- כל קבוצה בגודל  $d+1$  לא ניתנת לניתוח ע"י  $\mathcal{H}$

ככל שהמימד קטן כך המחלקה פחות מסובכת וקל יותר ללמוד. עבור  $\mathcal{H}$  סופית מתקיים  $VCdim(\mathcal{H}) \leq \log_2 |\mathcal{H}|$ . סיבוכיות הדגימות של למידת PAC על מחלקת המסווגים הבינאריים עם  $d = VCdim(\mathcal{H})$  היא:

$$C_1 \cdot \frac{d + \log(1/\delta)}{\epsilon} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \cdot \frac{d \cdot \log(1/\epsilon) + \log(1/\delta)}{\epsilon}$$

אם  $VCdim(\mathcal{H}) = \infty$  לא נוכל לבצע למידה. אם נתונות  $m < VCdim(\mathcal{H})$  דגימות ונרצה לחזות את הדגימה ה- $VCdim(\mathcal{H})$ ,  $m+1 \leq VCdim(\mathcal{H})$ , עדיין ניתן לנתץ את הקבוצה וקיימות שתי היפותאות שקולות שמסווגות באופן שונה את הדוגמא הנוספת.

## Support Vector Machine (SVM)

לא פרמטרי. נשאף להפרדה טובה - רחוקה כמה שיותר מכל הדגימות, margin מקסימלי אפשרי (מרחק הנקודות הקרובות ביותר אלין). נמקסם את  $J(w) = \frac{2}{\|w\|}$  תחת האילוץ  $y_i (w^T x_i + b) \geq 1$  לכל  $i$ . בשלב הסיווג נחזה את  $\hat{y} = \text{sign}(w^T x + b)$  כאשר:

$$w = \sum_{i=1}^n \alpha_i y_i x_i; \quad b = y_i - w^T x_i$$

דגימה  $x_i$  היא Support Vector (SV) אם  $\alpha_i > 0$ .

ב-soft SVM יהיה slack לכל דגימה - כמה היא יכולה לשבור את תנאי המרחק מה-hyperplane המפריד. נרצה למזער את ה-hinge loss:  $\sum_{i=1}^n \xi_i$ .  $\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$  תחת האילוץ  $\xi_i \geq \max\{0, 1 - y_i a^T x_i\}$ .  $C$  קטן - נאפשר slacks גדולים וכך שגיאת אימון גדולה יותר (בנוסף, יותר נקודות בתוך ה-margin). במילים אחרות, ההפסד הוא:

$$J(w, \xi_1, \dots, \xi_n) = C \sum_{i=1}^n \max\{0, 1 - y_i a^T x_i\} + \lambda \|w\|^2$$

עבור  $\lambda$  גדול נקבל  $\|w\|^2$  קטן, וכך margin גדול יותר.

כשהדאטה לא פריד לינארית נרצה לסבך את המודל ע"י גרעין  $K(x, y) = \varphi(x)^T \varphi(y)$ . נטיל את הנקודות למימד גבוה יותר באמצעות  $\varphi$  ונפריד לינארית במימד החדש:  $g(x) = w^T \varphi(x) + w_0$ . נעלה מימדים באופן ישיר ע"י הגרעין. צירוף לינארי של גרעינים הוא גרעין.

$$g(x) = \sum_{i=1}^n \alpha_i y_i K(x_i, x) + b$$

סיבוכיות המסווג תלויה במספר ה-SV שלו: רק עבורם  $\alpha > 0$  והמודל תלי רק בהם.

## Cross Validation

LOOCV: עבור על כל נקודה באימון, הסר אותה מהדאטה (אמנית), אמן את המודל על שאר הדאטה וחזה את הנקודה שנפלה. בסוף התהליך - חשב שגיאה ממוצעת. ב-k-fold בכל פעם נעיק חתיכה אחת (מתוך  $k$ ), נאמן על שאר הדאטה ונבדוק עבור מה שהעפנו. ככל שיש יותר מידע, נוכל ועדיף  $k$  קטן.