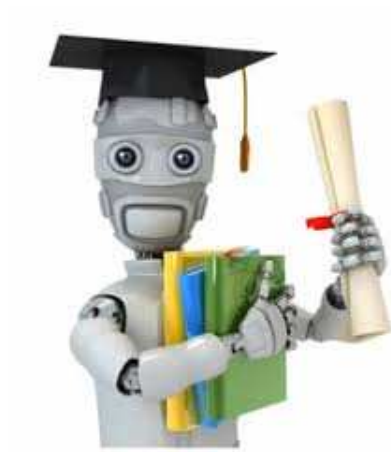


Introduction to Machine Learning (67577)

July 9, 2018



המסמך נכתב ע"י גלעד גרין, על בסיס הרצאות ותרגולי הקורס וספרי הקורס:

Understanding Machine Learning: From Theory to Algorithms

Elements of Statistical Learning

Reinforcement Learning: An Introduction

תוכן עניינים

6

I הקדמה

7

II ידע מוקדם

7	אלגברה לינארית	1
7	1.1 נורמה ומכפלה פנימית סטנדרטית	
8	1.2 פירוקי EVD, SVD	
9	1.3 שימוש SVD בדחיסה	
9	2 חשבון רב מימדי	
10	3 הסתברות	
10	3.1 אי שוויונות	
11	3.2 הסתברות רב מימדית	
11	4 $Convexity, Lipschitzness and Smoothness$	
11	4.1 $Convexity$	
11	4.1.1 קבוצות קמורות	
12	4.1.2 פונקציות קמורות	
14	4.2 ליפשיציות	
14	4.3 חלקות	

15

III תיאוריה כללית

15	5 יסודות בלמידה	
15	5.1 סוגי למידה	
15	5.2 המודל של הלמידה הסטטיסטית	
16	5.3 $ERM, Realizability and i.i.d$	
16	5.3.1 $Loss functions$	
17	5.3.2 עיקרון ERM	
17	5.3.3 $The realizability assumption$	
17	5.3.4 $Independently and Identically Distributed - i.i.d$	
17	5.4 $No Free Lunch$	
17	5.5 $Classification Errors$	
17	5.5.1 $Accuracy$	

18 סוגי שגיאות	5.5.2	
18 Precision/Recall	5.5.3	
19 ROC עקומת	5.6	
20 PAC and Agnostic PAC		6
20 PAC Model	6.1	
20 Sample Complexity	6.2	
21 Agnostic PAC	6.3	
21 Data Generating Distribution	6.3.1	
21 True error/Empirical Risk	6.3.2	
22 דוגמה - Coin Prediction	6.4	
23 VC – Dimension		7
24 Uniform Convergence		8
25 המשפט היסודי של הלמידה הסטטיסטית		9
26 Bias – Complexity Tradeoff הבנת		10
26 Error Decomposition	10.1	
27 Tradeoff	10.2	
27 דוגמה עבור פולינומים	10.3	

29 תיאוריה לאלגוריתמים	IV	
29 מסווגים לינאריים		11
29 חצאי מרחב (Halfspaces)		11.1
29 ERM עבור halfspaces באמצעות Perceptron	11.1.1	
30 ERM עבור halfspaces באמצעות תכנון לינארי	11.1.2	
31 VCdim של halfspaces	11.1.3	
31 רגרסיה לינארית		11.2
31 פתרון הבעיה	11.2.1	
34 קיומו של רעש	11.2.2	
35 Numerically Stable	11.2.3	
35 Bias – Variance Decomposition	11.2.4	
36 רגרסיה לוגיסטית		11.3
36 אינטואיציה	11.3.1	
37 הגדרה פורמלית	11.3.2	
38 multiclassification מקרה	11.3.3	
38 Bayes Classifiers		12
38 Bayes Optimal	12.1	
39 Naive Bayes	12.2	

39	<i>Linear Discriminant Analysis</i>	12.3
41	<i>Quadratic Discriminant Analysis</i>	12.4
42	<i>Boosting</i>	13
42	<i>Weak Learnability</i>	13.1
43	<i>Boosting Confidence</i>	13.2
44	<i>Boosting Confidence בתוחלת</i>	13.2.1
44	<i>AdaBoost באמצעות Boosting Accuracy</i>	13.3
47	<i>Bias – complexity tradeoff</i>	13.3.1
48	<i>Model Selection and Validation</i>	14
48	<i>Validation and k – Cross Validation</i>	14.1
49	<i>k – fold בחירת ערך ה-k עבור</i>	14.1.1
50	<i>k איך לא בוחרים</i>	14.1.2
51	<i>Train – Validation – Test split</i>	14.2
52	<i>Regularization</i>	15
52	<i>RLM – Regularized Loss Minimization עיקרון</i>	15.1
52	<i>Best Subset Regression</i>	15.1.1
53	<i>Ridge Regression</i>	15.1.2
55	<i>L₁ Lasso</i>	15.1.3
56	<i>השוואה בין פונקציות הרגולריזציה</i>	15.2
56	<i>Hard/Soft Thresholding</i>	15.2.1
57	<i>sparsity כדורי יחידה של הנורמות ו-</i>	15.2.2
58	<i>Logistic Regression עבור Lasso</i>	15.3
59	<i>Bootstrapping and Bagging</i>	16
59	<i>Bootstrapping</i>	16.1
59	<i>Bagging</i>	16.2
60	<i>Feature Selection</i>	17
62	<i>SVM</i>	18
62	<i>Margin and Hard – SVM</i>	18.1
64	<i>המקרה ההומוגני</i>	18.1.1
64	<i>Hard – SVM של Sample Complexity</i>	18.1.2
65	<i>Soft – SVM</i>	18.2
66	<i>SVM vs. Boosting</i>	18.3
66	<i>Kernels</i>	19
66	<i>Embedding into Feature Space</i>	19.1
67	<i>Kernel Trick</i>	19.2
69	<i>Decision Trees</i>	20
70	<i>Bias – Complexity tradeoff</i>	20.1

70	אלגוריתמים	20.2
70	<i>CART – Classification and Regression Trees</i>	20.2.1
71	<i>Pruning</i> (גזום)	20.3
72	<i>Bias variance</i>	20.3.1
72	שימוש ב- <i>Bagging</i> ו- <i>Bootstrapping</i>	20.4
73	<i>Random Forest</i>	20.5
74	<i>k – Nearest Neighbors</i>	21
74	האלגוריתם	21.1
74	<i>Bias – complexity tradeoff</i>	21.2
75	<i>Bellman's "Curse of dimensionality"</i>	21.3
75	<i>Convex Learning Problems</i>	22
75	בעיות למידה	22.1
76	פונקציית איבוד <i>Surrogate</i>	22.2
77	<i>Stochastic Gradient Descent</i>	23
77	<i>Gradient Descent</i>	23.1
78	<i>Sub Gradient Descent</i>	23.2
80	אלגוריתם <i>Stochastic Gradient Descent (SGD)</i>	23.3
80	<i>SGD</i> עבור מינימום L_D	23.3.1
81	<i>SGD</i> עבור מינימום L_S (ERM)	23.3.2
83	<i>Deep Learning</i>	24
84	<i>Computation Graph and Backpropagating</i>	24.1
87	<i>Expressiveness and Sample Complexity</i>	24.2
87	<i>Tricks</i>	24.3
88	<i>Convolutional Networks</i>	24.4
89	<i>Reinforcement Learning</i>	25
89	הקדמה	25.1
91	<i>Finite Markov Decision Processes (MDP)</i>	25.2
92	<i>The Markovian Assumption</i>	25.2.1
92	<i>Policies and Value functions</i>	25.2.2
94	פוליסות ופונקציות ערך אופטימליות	25.2.3
95	אלגוריתמים ל- <i>MDP</i>	25.3
96	<i>Policy Evaluation</i>	25.3.1
97	<i>Value Iteration</i>	25.3.2
98	<i>Learning</i>	25.3.3
100	<i>Dimensionality Reduction</i>	26
100	<i>Linear Dimension Reduction and PCA</i>	26.1
102	<i>Clusters – k – means</i>	27

חלק I

הקדמה

תהליך הלמידה עוסק בהפיכת ניסיון לידע או מומחיות. בהינתן עולם בעיה (*domain*) נרצה אלגוריתם (*learner*) אשר ילמד סט מידע (*trainset*) לספק פרדיקציה טובה עבור פריטים חדשים באותו עולם תוכן (*testset*). כאשר ניגש לבעיית למידה נרצה ללמוד את ה-*data*. נרצה להבין איך ה-*domain* מתנהג, מהן התכונות (*features*) שנרצה ללמוד בעזרתן. מתוך כך נרצה להעריך איזה מודל מתאים ל-*domain* (*model selection*) ולבסוף להעריך את טיב המודל (*model – assesment*) על מידע שכלל עוד לא ראינו (*validation set*). במהלך הקורס נפגוש בכלים מתמטיים רבים. כאשר מדברים על כיצד המידע מתנהג למעשה מדובר בשאלה הסתברותית. כיצד המידע מתפלג על גבי ה-*features* השונים והתיוגים (*labels*) השונים. כאשר נרצה ללמוד *domain* מסויים נרצה מסווג שיטעה (*loss*) כמה שפחות והרי שמכאן כי זאת בעיית מינימיזציה אשר בהתאם למודלים השונים נפגוש חשבון אינפיטיסימלי ואלגברה לינארית.

מטרת הקורס ללמד איך לא להיות "machine learning monkey". כלומר נרצה להבין את המודלים השונים ולהחליט איזה מודל מתאים לנו במקרה הספציפי של ה-*domain* עליו אנו עובדים. "להבין" בקורס זה משלב את ההבנה הפורמלית של משפטים וטענות, ואת פיתוח האינטואיציה עבור התהליכים השונים.

על אף שסיכום זה מבוסס על החומר שהועבר בהרצאות ובתרגולים, סדר הנושאים מנסה להתאים בין הסדר בספר הלימוד לבין סדר ההרצאות.

חלק II

ידע מוקדם

פרק זה עוסק ביסודות המתמטיים הנחוצים לקורס ומכיל את ההגדרות שהועברו בתרגולים בנושאי אלגברה לינארית, חשבון אינפיניטסימלי במרחבים מממד גבוהה והסתברות.

1 אלגברה לינארית

1.1 נורמה ומכפלה פנימית סטנדרטית

הגדרה 1.1 נורמה הינה פונקציה $\|\cdot\| : \mathbb{R}^m \rightarrow \mathbb{R}$ המקיימת את התכונות הבאות לכל $a \in \mathbb{R}$, $u, v \in \mathbb{R}^m$:

$$1. \text{ חיוביות - } \|v\| \geq 0 \text{ ואם } \|v\| = 0 \text{ אז } v = 0.$$

$$2. \text{ הומוגניות חיובית - } \|av\| = |a| \cdot \|v\|$$

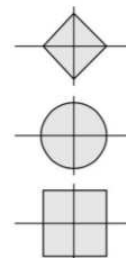
$$3. \text{ אי שיוויון המשולש - } \|u + v\| \leq \|u\| + \|v\|$$

הגדרה 1.2 כדור היחידה הינו $B = \{x \mid \|x\| \leq 1\}$. כאשר בהתאם לסוג הנורמה התוצאה שמתקבלת היא שונה:

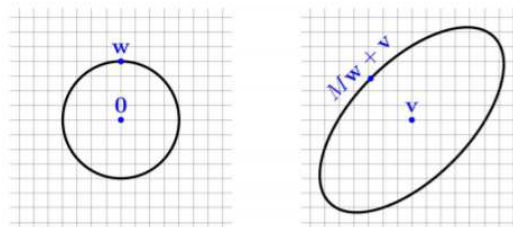
$$\|x\|_1 = \sum_{i=1}^m |x_i|,$$

$$\|x\|_2 = \left(\sum_{i=1}^m |x_i|^2 \right)^{1/2} = \sqrt{x^* x},$$

$$\|x\|_\infty = \max_{1 \leq i \leq m} |x_i|,$$



הגדרה 1.3 אליפסואיד הינו סט נקודות $A = \{y \mid y = Mx + v \wedge \|x\|^2 = 1\}$ כאשר M מטריצה הפיכה.



טענה 1.4 תהי U מטריצה אורתוגונלית אזי האליפסואיד המוגדר על ידי A שווה לאליפסואיד המוגדל על ידי AU .

הגדרה 1.5 מכפלה פנימית הינה פונקציה $\langle \cdot, \cdot \rangle : \mathbb{R}^m, \mathbb{R}^m \rightarrow \mathbb{R}$ המקיימת את התכונות הבאות: לכל $a \in \mathbb{R}$, $x, y, z \in \mathbb{R}^m$

1. סימטריות - $\langle x, y \rangle = \langle y, x \rangle$

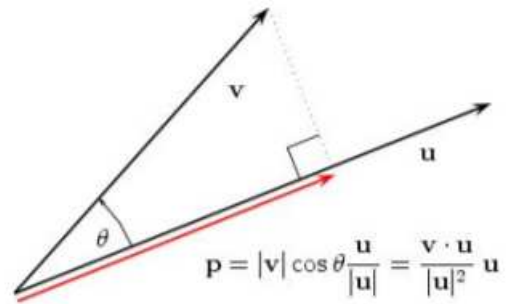
2. לינאריות - $\langle ax + z, y \rangle = a \langle x, y \rangle + \langle z, y \rangle$

3. חיוביות: $\langle x, x \rangle \geq 0$ ומתקיים $\langle x, x \rangle = 0 \Leftrightarrow x = 0$

הגדרה 1.6 המכפלה הפנימית הסטנדרטית הינה: $\langle x, y \rangle = x^T \cdot y = \sum_{i=1}^m x_i \cdot y_i$.

טענה 1.7 עבור זוג וקטורים x, y מתקיים כי המכפלה הסקלרית שלהם היא $\langle x, y \rangle = \|x\| \cdot \|y\| \cdot \cos \theta$ כאשר θ הזווית שבין הוקטורים.

הגדרה 1.8 ההטלה של וקטור v על וקטור u הינה הוקטור p עם אורך $\|v\| \cos \theta$ בכיוון הוקטור u :



יהי $w \in \mathbb{R}^d$ אזי ה- $halfspace$ הינו הקבוצה $\{x | \langle w, x \rangle \geq b\}$ וה- $hyperplane$ (על מישור) הינו הקבוצה $\{x | \langle w, x \rangle = b\}$.

1.2 פירוקי EVD, SVD

משפט 1.9 פירוק EVD - תהי $B \in \mathbb{R}^{m \times m}$ סימטרית אזי קיימת U אורתוגונלית ו- λ אלכסונית כך ש- $B = U \lambda U^T$ ולכל $u_i \in U$ וקטור עצמי של B עם ערך עצמי $\lambda_{i,i}$.

משפט 1.10 פירוק SVD - תהי $A \in \mathbb{R}^{d \times m}$ אזי קיימת $V \in \mathbb{R}^{m \times m}, U \in \mathbb{R}^{d \times d}$ אורתוגונליות ו- $\Sigma \in \mathbb{R}^{d \times m}$ אלכסונית כך ש- $A = U \Sigma V^T$. בנוסף מתקיים שאם $rank(A) = r$ אזי לכל $i \leq r$ הוקטורים v_i, u_i הינם הוקטורים הסינגולריים של A עם ערך סינגולרי $\sigma_i \in \Sigma_{i,i}$ ולכל $i > r$ אז $A^T u_i = 0, A v_i = 0$.

טענה 1.11 תכונות של SVD :

1. σ_i^2 הם ערכים סינגולריים של המטריצה XX^T .

2. u_i הם וקטורים סינגולריים ימניים המתאימים של XX^T .

3. v_i הם וקטורים סינגולריים שמאליים המתאימים של $X^T X$.

1.3 שימוש SVD בדחיסה

נניח כי $X \in \mathbb{R}^{d \times m}$ ורוצים לשערך אותה על ידי מטריצה מדרגה $k \ll d, m$. נרצה למצוא $\tilde{X} \in \mathbb{R}^{d \times m}$ שדרגתה k הממוזערת את:

$$\tilde{r} = \text{rank}(X - \tilde{X}), \|X - \tilde{X}\|_F = \sum_{i,j} (X_{i,j} - \tilde{X}_{i,j})^2 = \sum_{i=1}^{\tilde{r}} \sigma_i(X - \tilde{X})$$

כבעיית אופטימיזציה נוכל לכתוב בתור $\argmin_{\tilde{X} : \text{rank}(\tilde{X}) \leq k} \|X - \tilde{X}\|_F$.

משפט 1.12 תהי $X \in \mathbb{R}^{d \times m}$ עם SVD $X = U \Sigma V^T$. נגדיר $\tilde{\Sigma} \in \mathbb{R}^{d \times m}$ אלכסונית בתור

$$\tilde{\Sigma}_{i,i} = \begin{cases} \Sigma_{i,i} & i < k \\ 0 & i \geq k \end{cases}$$

אזי הפתרון לבעיית הדחיסה הוא $\tilde{X} = U \tilde{\Sigma} V^T$.

2 חשבון רב מימדי

הגדרה 2.1 הנגזרת החלקית של הפונקציה $f : \mathbb{R}^n \rightarrow \mathbb{R}$ בנקודה $x \in \mathbb{R}^n$ ביחס ל- x_i היא:

$$\frac{\partial f(x)}{\partial x_i} = \lim_{a \rightarrow 0} \frac{f(x + a \cdot e_i) - f(x)}{a}$$

הגדרה 2.2 הגרדיאנט של פונקציה $f : \mathbb{R}^n \rightarrow \mathbb{R}$ בנקודה $x \in \mathbb{R}^n$ הוא:

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$$

הערה 2.3 הגרדיאנט מספק תנאי הכרחי לנקודות קיצון של פונקציה. כלומר אם x נקודת קיצון והגרדיאנט בנקודה קיים אזי $\nabla f(x) = 0$.

הגדרה 2.4 תהי $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ו- $x_0 \in \mathbb{R}^n$ אזי הקירוב הליניארי של f סביב x_0 הוא $f(x_0) + \nabla f(x_0) \cdot (x - x_0)$. כלומר לכל $x \in \mathbb{R}^n$ מתקיים כי: $f(x) \approx f(x_0) + \nabla f(x_0) \cdot (x - x_0)$.

הגדרה 2.5 קו הגובה של $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ביחס ל- $c \in \mathbb{R}$ הוא הקבוצה $\{x | f(x) = c\}$.

הגדרה 2.6 עבור פונקציה $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$ היעקוביאן של הפונקציה הינה מטריצה $x \times d$ של כל הנגזרות החלקיות:

$$J_x(f) = \begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1} & \dots & \frac{\partial f_1(x)}{\partial x_d} \\ \vdots & & \vdots \\ \frac{\partial f_m(x)}{\partial x_1} & \dots & \frac{\partial f_m(x)}{\partial x_d} \end{bmatrix}$$

דוגמה: נגדיר כדור ϵ סביב נקודה x_0 : $B_\epsilon(x_0) = \{x_0 + x | \|x\|_2 = \epsilon\}$. על כן מהו הוקטור שממקסם/ממזער את הקירוב הלינארי של f סביב x_0 על פני הכדור?

פתרון: נראה כי מתקיים:

$$f(x + x_0) \approx f(x_0) + \nabla f(x_0)^T x = \underbrace{\|\nabla f(x_0)\|_2}_{\text{doesn't depend on } x} \cdot \underbrace{\|x\|_2}_{\epsilon} \cdot \cos\theta$$

ולכן מזעור/מקסום הקירוב הלינארי שקול למיקסום/מזעור $\cos\theta$ כאשר זוהי הזווית בין x לגרדיאנט של f בנקודה x_0 .
על כן נקבל כי עבור:

$$\begin{aligned} \bullet \quad \theta = 0^\circ & \quad x = \epsilon \frac{\nabla f(x_0)}{\|\nabla f(x_0)\|_2} \quad \text{הוקטור הממקסם} \\ \bullet \quad \theta = 180^\circ & \quad x = -\epsilon \frac{\nabla f(x_0)}{\|\nabla f(x_0)\|_2} \quad \text{הוקטור הממזער} \end{aligned}$$

3 הסתברות

3.1 אי שוויונות

משפט 3.1 אי שוויון מרקוב - יהי X משתנה מקרי (מ"מ) אי שלילי בעל תוחלת סופית אזי לכל $0 < a \in \mathbb{R}$ מתקיים כי:

$$\mathbb{P}[X \geq a] \leq \frac{\mathbb{E}[X]}{a}$$

משפט 3.2 אי שוויון צ'בישב - יהי X מ"מ בעל תוחלת ושוונות סופיות אזי לכל $0 < a \in \mathbb{R}$ מתקיים כי:

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq a] \leq \frac{\text{Var}(X)}{a^2}$$

משפט 3.3 אי שיויון הופדינג - יהיו X_1, \dots, X_m מ"מ בלתי תלויים (ב"ת) וחסומים כך שלכל $i \in [m]$ מתקיים $a_i \leq X_i \leq b_i$. יהי $\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i$ אזי

$$\mathbb{P}[|\bar{X} - \mathbb{E}[\bar{X}]| \geq \epsilon] \leq 2 \exp\left(\frac{-2m^2\epsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}\right)$$

3.2 הסתברות רב מימדית

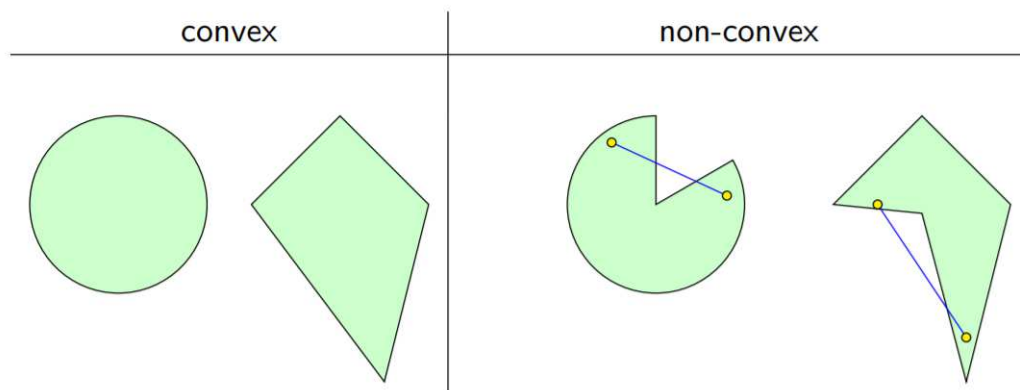
להשלים

4 Convexity, Lipschitzness and Smoothness

4.1 Convexity

4.1.1 קבוצות קמורות

הגדרה 4.1 יהי V מרחב וקטורי. $C \subseteq B$ נקראת קמורה אם לכל $u, v \in C$ ו- $\alpha \in [0, 1]$ מתקיים $\alpha u + (1 - \alpha)v \in C$



הגדרה 4.2 עבור קבוצת נקודות $\{v_i \mid i \in I\}$ הנקודה $\sum_{i \in I} \alpha_i v_i$, כאשר המקדמים $\{\alpha_i > 0 \mid i \in I\}$ חיוביים ונסכמים ל-1, נקראת צירוף קמור של ה- v_i .

דוגמה: כדורי יחידה הינם קבוצה קמורה.

הוכחה: יהי $B = \{v \in V \mid \|v\| = 1\}$. יהי $u, v \in B$ ו- $\alpha \in [0, 1]$ אז:

$$\begin{aligned} \|\alpha u + (1 - \alpha)v\| &\leq \|\alpha u\| + \|(1 - \alpha)v\| \\ &= \underbrace{\alpha \|u\| + (1 - \alpha)\|v\|}_{\downarrow} \leq \alpha + 1 - \alpha = 1 \\ &\downarrow \\ \alpha u + (1 - \alpha)v &\in B \end{aligned}$$

משפט 4.3 כל הקבוצות הבאות הן קמורות:

1. חיתוך קבוצות קמורות הוא קמור. $C = \bigcap_{i \in I} C_i$ עבור $\{C_i \mid i \in I\}$ אוסף קבוצות קמורות.

2. סכום וקטורי $C_1 + C_2 = \{c_1 + c_2 \mid c_1 \in C_1, c_2 \in C_2\}$ עבור C_1, C_2 קמורות.

3. הקבוצה $\lambda C = \{\lambda c \mid c \in C\}$ עבור C קמורה.

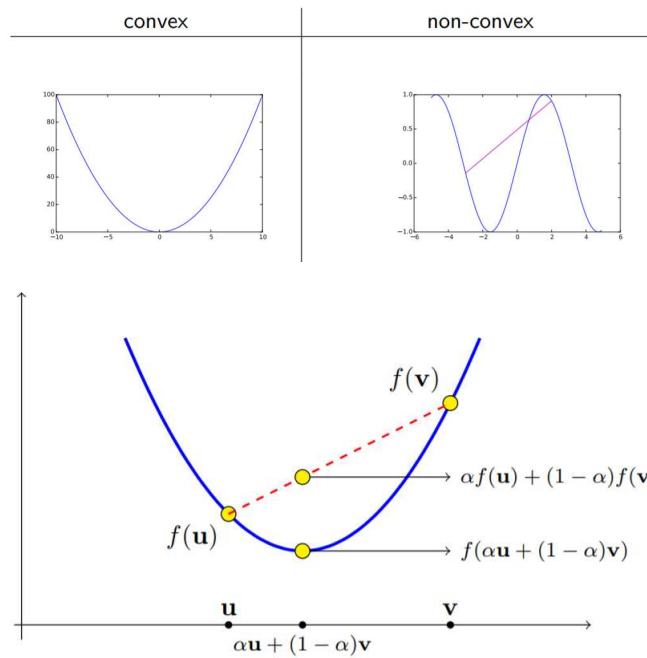
הוכחה: (מקרה 2) יהיו $u, v \in C_1 + C_2$. נראה שלכל $\alpha \in [0, 1]$ מתקיים $\alpha u + (1 - \alpha)v \in C_1 + C_2$. נראה כי מתקיים $u = a_1 + a_2, v = b_1 + b_2$ כך ש- $a_1, b_1 \in C_1$ ו- $a_2, b_2 \in C_2$. כעת:

$$\begin{aligned} \alpha u + (1 - \alpha)v &= \alpha(a_1 + a_2) + (1 - \alpha)(b_1 + b_1) \\ &= \underbrace{\alpha a_1 + (1 - \alpha)b_1}_{\in C_1} + \underbrace{\alpha a_2 + (1 - \alpha)b_2}_{\in C_2} \\ &\downarrow \\ \alpha u + (1 - \alpha)v &\in C_1 + C_2 \end{aligned}$$

4.1.2 פונקציות קמורות

הגדרה 4.4 יהי מרחב וקטור V ו- C קבוצה קמורה. פונקציה $f : C \rightarrow \mathbb{R}$ תקרא קמורה אם לכל $u, v \in C$ ו- $\lambda \in [0, 1]$ מתקיים:

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



הגדרה 4.5 ה- $epigraph$ של פונקציה היא הקבוצה $epigraph(f) = \{(x, \beta) \mid f(x) \leq \beta\}$

משפט 4.6 פונקציה f הינה קמורה $\iff ephgraph(f)$ קבוצה קמורה

דוגמה: נורמה $\|\cdot\|$ היא פונקציה קמורה. לכל $u, v \in \mathbb{R}^d$ ול- $\lambda \in [0, 1]$ מתקיים:

$$\|\lambda u + (1 - \alpha) \lambda\| \stackrel{tri. inequality}{\leq} \|\lambda u\| + \|(1 - \lambda) v\| = \lambda \|u\| + (1 - \lambda) \|v\|$$

דוגמה: ההעתקה האפינית $f(u) = \langle w, u \rangle + b$ הינה פונקציה קמורה.

משפט 4.7 הפונקציה $g(u)$ הינה פונקציה קמורה בכל המקרים הבאים:

1. סכום ממושקל חיובית של פונקציות קמורות: $g(u) = \sum_{i \in I} \gamma_i f_i(u)$ עבור $\{\gamma_i \mid i \in I\}$ חיוביות ו- $\{f_i \mid i \in I\}$ קמורות.

2. הרכבה של פונקציה קמורה על אפינית: $g(u) = f(Au + b)$ עבור f קמורה ו- $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$.

3. sup או max של אוסף פונקציות קמורות: $g(u) = \sup_{i \in I} \{f_i(u)\}$ עבור $\{f_i \mid i \in I\}$ קמורות.

הוכחה: (מקרה 3) נשים לב כי $\forall i \in I \quad f_i(u) \leq t \iff (u, t) \in \text{epigraph}(f_i)$ ומכאן כי $\text{epigraph}(g) = \bigcap_{i \in I} \text{epigraph}(f_i)$. היות ו- f_i קמורות אזי חיתוך קמורות הוא קמור ולכן גם g קמורה.

4.2 ליפשיציות

הגדרה 4.8 תהי $C \subset \mathbb{R}^d$. פונקציה $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$ הינה ρ -ליפשיצית מעל C אם $\forall w_1, w_2 \in C$ מתקיים $\|f(w_1) - f(w_2)\| \leq \rho \|w_1 - w_2\|$

באופן אינטואיטיבי פונקציה ליפשיצית מוגבלת במהירות ההשתנות שלה

דוגמה: פונקציה לינארית $f : \mathbb{R}^d \rightarrow \mathbb{R}$ המוגדרת על ידי $f(w) = \langle v, w \rangle + b$ עבור $v \in \mathbb{R}^d$ הינה $\|v\|$ -ליפשיצית. בעזרת אי שיוויון קושי-שוורץ נקבל:

$$|f(w_1) - f(w_2)| = |\langle v, w_1 - w_2 \rangle| \leq \|v\| \cdot \|w_1 - w_2\|$$

טענה 4.9 תהי $f(x) = g_1(g_2(x))$ כאשר g_1 הינה ρ_1 -ליפשיצית ו- g_2 הינה ρ_2 -ליפשיצית אזי f הינה $\rho_1 \rho_2$ -ליפשיצית

4.3 חלקות

הגדרה 4.10 פונקציה דיפרנציאבילית $f : \mathbb{R}^d \rightarrow \mathbb{R}$ הינה β -חלקה אם הגרדיאנט שלה הינו β -ליפשיצי. כלומר $\forall v, u \quad \|\nabla f(v) - \nabla f(u)\| \leq \beta \|v - u\|$

טענה 4.11 תהי $f(w) = g(\langle w, x \rangle + b)$ כאשר g פונקציה β -חלקה ו- $x \in \mathbb{R}^d$ אזי f הינה $\beta \|x\|^2$ -חלקה. (הוכחה עמוד 163 בספר הלימוד)

טענה 4.12 תהי פונקציה קמורה ו- β -חלקה אזי $\|\nabla f(w)\|^2 \leq 2\beta f(w)$. פונקציה המקיימת תכונה זו נקראת גם פונקציה *self-bounded*.

חלק III

תיאוריה כללית

5 יסודות בלמידה

5.1 סוגי למידה

- Batch vs. Online - בסוג למידה *batch* נקבל *training set* מתוויג ו- *test set* שאיננו מתוויג. נלמד על סט האימון ונספק *estimation* עבור סט הבדיקה. בלמידה *online* בכל פעם נקבל *sample*, נעדכן את המודל בהתאם ל-*sample* החדש ונחזיר פרדיקציה.
- Supervised vs. Unsupervised - כאשר ה-*trainset* מתוויג בתיוג "אמת" נכנה את הלמידה *supervised*. עבור *trainset* שאיננו מתוויג נכנה את הלמידה *unsupervised*.

5.2 המודל של הלמידה הסטטיסטית

הקלט שניתן ללומד:

1. *Domain set* הינו סט ארביטררי \mathcal{X} של מידע המכיל דוגמיות המכונות *samples*. כל $x \in \mathcal{X}$ מייצג *feature vector*.
2. *Label Set* הינו סט התיוגים האפשריים \mathcal{Y} עבור המידע ב- \mathcal{X} .
3. *Training Set* הינו אוסף $S = \{(x_i, y_i)\}_{i=1}^m \subseteq \mathcal{X} \times \mathcal{Y}$ כאשר $\forall i \in [m] x_i \in \mathcal{X}, y_i \in \mathcal{Y}$.

בתהליך הלמידה:

1. נרצה שהלומד יתאמן על המידה ויחזיר *Prediction Rule* - פונקציה מסט המידע לסט התיוגים. $h : \mathcal{X} \rightarrow \mathcal{Y}$. נסמן את פלט האלגוריתם $\mathcal{A}(S)$.
2. בהינתן דומיין \mathcal{X} וסט תיוגים \mathcal{Y} מחלקת היפותזות \mathcal{H} הינה קבוצת כל הפונקציות $\mathcal{H} = \{h : \mathcal{X} \rightarrow \mathcal{Y}\}$.

כאשר \mathcal{D} ותויג נכונה על ידי פונקציה כלשהי f . עבור חוק הפרדיקציה שהחזיר הלומד נרצה למדוד את הטעות של חוק הפרדיקציה h . נגדיר את הטעות של מסווג, ה-*true error* או ה-*generalization error* על ידי:

$$L_{\mathcal{D},f}(h) \stackrel{\text{def.}}{=} \mathbb{P}_{x \sim \mathcal{D}}[h(x) \neq f(x)] \stackrel{\text{def.}}{=} \mathcal{D}(\{x : h(x) \neq f(x)\})$$

בהמשך נראה כי הגדרה זו משתנה מעט. היות והלומד איננו יודע את \mathcal{D}, f נוכל רק לבדוק את ה-*training error*. נקרא גם *empirical error* או *empirical risk*:

$$L_S(h) \stackrel{\text{def.}}{=} \frac{|\{i \in [m] : h(x_i) \neq y_i\}|}{m}$$

נשים לב שבהגדרה זו אנו מניחים שמשקל כל טעות שווה - התפלגות אחידה.

5.3 ERM, Realizability and i.i.d

כאשר תיארונו את הרעיון המנחה בלמידה רצינו כי הלומד ידע לספק פרדיקציה טובה עבור המידע. ההגדרה של "טובה" איננה חד משמעית אולם לרוב מכילה גם אינפורמציה עבור ה"איבוד" של חוק הפרדיקציה שנחזיר. כלומר החזרת אמת-מידה לטעויות הפרדיקציה שביצע האלגוריתם.

5.3.1 Loss functions

קיימות פונקציות איבוד רבות בעלות תכונות שונות והמשמשות אלגוריתמים שונים:

- 0-1 Loss - פונקציית איבוד אשר נותנת 1 עבור טעות ו-0 להצלחה. בהינתן סט מתוייג $S = \{(x_i, y_i)\}_{i=1}^m$ נגדיר:

$$L_S^{0-1}(h) = \sum_{i=1}^m \mathbb{1}_{h(x_i) \neq y_i}$$

- Absolute Value Loss - פונקציית איבוד הסוכמת את מרחקי הפרדיקציות מערכי האמת:

$$L_{S,f} = \sum_{(x,y) \in S} |y - f(x)|$$

- Least Squares Loss - פונקציית איבוד הסוכמת את ריבועי הפרשי המרחקים בין הפרדיקציות לערכי האמת:

$$L_{S,f} = \sum_{(x,y) \in S} (y - f(x))^2$$

5.3.2 עיקרון ERM

אלגוריתם למידה המממש את עיקרון ה-ERM Empirical Risk Minimization – מחזיר היפותזה הממזער את שגיאת האימון. עבור $f: \mathcal{X} \rightarrow \mathcal{Y}$ פונקציית תיוג המתייגת נכונה את \mathcal{X} :

$$ERM_{\mathcal{H}}(S) \in \underset{h \in \mathcal{H}}{\operatorname{argmin}} L_S(h) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{m} |\{x_i \in S \mid h(x_i) \neq f(x_i)\}|$$

5.3.3 The realizability assumption

הגדרה 5.1 עבור מחלקת היפותזות \mathcal{H} הנחת ה-*realizability* אומרת כי קיימת $h^* \in \mathcal{H}$ כך שמתקיים $L_{\mathcal{D},f}(h^*) = 0$. מכאן גם כי לכל $S \subseteq \mathcal{X}$ שנדגם בהתאם ל- \mathcal{D} מתקיים $L_S(h^*) = 0$.

5.3.4 Independently and Identically Distributed – i.i.d

כאשר נבנה סט אימון מתוך \mathcal{X} אשר בעל פילוג \mathcal{D} נבחר את ה-*samples* באופן מקרי ונניח כי הם בלתי תלויים אחד בשני. נסמן זאת בתור $S \sim \mathcal{D}^m$ אשר אומר למעשה כי x_1, \dots, x_m הם m "ב"ת וש"ה מתוך הפילוג \mathcal{D} .

5.4 No Free Lunch

משפט 5.2 יהי אלגוריתם למידה \mathcal{A} המתייג בינארית וביחס לפונקציית האיבוד $0-1$ מעל דומיין \mathcal{X} . יהי $m < \frac{|\mathcal{X}|}{2}$ מספר ה-*samples* בסט האימון S . אזי קיימת התפלגות \mathcal{D} מעל $\mathcal{X} \times \{0, 1\}$ כך ש:

$$1. \quad L_{\mathcal{D}}(f) = 0 \text{ עבורה } f: \mathcal{X} \rightarrow \{0, 1\} \text{ קיימת}$$

$$2. \quad \text{בהסתברות לפחות } \frac{1}{7} \text{ מעל בחירה של } S \sim \mathcal{D}^m \text{ נקבל כי } L_{\mathcal{D}}(\mathcal{A}(S)) \geq \frac{1}{8}$$

כלומר, לכל *learner* תהיה משימת למידה עליה הוא יכשל.

הערה 5.3 עבור אלגוריתם המתייג אקראי נקבל כי האיבוד שווה ל- $\frac{1}{2}$. מכאן כי המשפט אומר כי לא נוכל לעשות טוב יותר מניחוש אקראי.

5.5 Classification Errors

5.5.1 Accuracy

הדרך הפשוטה ביותר לבדוק את טיב הסיווג הוא על ידי מדד זה. עבורו פשוט נסכום את סך הטעויות בסיווג שביצע המסווג.

$$\sum_{i=1}^m \mathbb{1}_{y_i \neq \hat{f}(x_i)}$$

5.5.2 סוגי שגיאות

כאשר אנו מתעסקים עם בעיות *classification* מתקבלות שתי סוגי שגיאות: אמרנו a כאשר למעשה היה b או אמרנו b כאשר למעשה היה a . לשתי טעויות אלה אופי שונה בהתאם לסוג הבעיה שאנו מנסים לפתור ועלינו להחליט איזה סיווג הינו בעל אפקט ואיזה איננו בעל אפקט. נניח כי אנו מסווגים *data* מסויים ל"מעניין" ו"לא מעניין" ונניח כי איננו יכולים להסכים לטעויות רבות שבו אמרו עם דוגמאות שהן "לא מעניינות" כאשר הן למעשה כן. על כן נסמן תיוג $y = -1$ בתור "מעניין" ו- $y = 1$ בתור "לא מעניין" אזי עבור \hat{y} תוצאת הסיווג על ידי המסווג נקבל:

	$y_i = -1$	$y_i = 1$
$\hat{y}_i = -1$.	Type-II error
$\hat{y}_i = 1$	Type-I error	.

על כן

- *Type - I error* - אמרנו "יש אפקט" כאשר "אין אפקט". טעות זו הינה פחות חשובה. בדוגמה למעשה אמרו כי הדוגמה "מעניינת" כאשר למעשה היא הייתה "לא מעניינת" והגדרנו כי נרצה כמה שפחות לטעות על המקרים שבהם טענו כי הדוגמה "לא מעניינת"
- *Type - II error* - אמרנו "אין אפקט" כאשר למעשה "יש אפקט". טעות זו יותר משמעותית. בדוגמה למעשה אמרנו כי דוגמה "לא מעניינת" כאשר היא הייתה "מעניינת" - סוג הטעות שרצינו להפחית.

כעת כשנסתכל שוב על הטבלה נוכל לומר כי:

	$y_i = -1$	$y_i = 1$
$\hat{y}_i = -1$	true negative	false negative
$\hat{y}_i = 1$	false positive	true positive

5.5.3 Precision/Recall

נסמן P את מספר הדוגמאות ה-*positive* ו- N מספר הדוגמאות ה-*negative* אזי:

• $(FP + FN) / (P + N)$ - *Error Rate*

• $(TP + TN) / (P + N)$ - *Accuracy*

• $TP / (TP + FP)$ - *Precision*

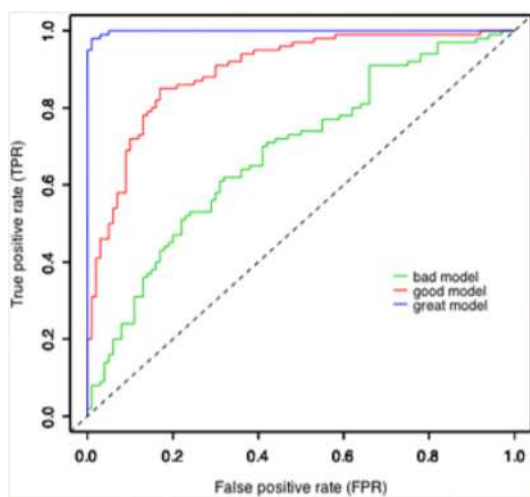
• TP / P - *Recall (True positive rate)*

• TN/N - *Specificity*

• FP/N - *False positive rate*

5.6 עקומת ROC

עקומה זו משמשת מניתוח ה-*tradeoff* בין ה-*sensitivity* (ה-*true positive rate, recall*) של המודל לבין ה-*specificity* של המודל.



PAC and Agnostic PAC 6

בלמידה הסטטיסטית עבור \mathcal{X} domain, קבוצת תגיות \mathcal{Y} ומחלקת היפותזות \mathcal{H} ישנו פילוג \mathcal{D} על פני \mathcal{X} אשר איננו ידוע ל- learner . ממשפט ה-*No Free Lunch* לא נוכל לצפות לקבל כפלט מהאלגוריתם $A(S)$ פרדיקציה אשר ללא איבוד על הפילוג הלא ידוע \mathcal{D} . על כן נרצה להיות מסוגלים להחזיר פרדיקציה שמצליחה "טוב מספיק" בהסתברות "גבוהה מספיק". כעת נבטא באופן פורמלי רצון זה.

PAC Model 6.1

הגדרה 6.1 מחלקת היפותזות \mathcal{H} הינה *PAC learnable* אם קיים אלגוריתם A ופונקציה $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ עם התכונות הבאות: לכל $\epsilon, \delta \in (0, 1)$, לכל פילוג \mathcal{D} מעל \mathcal{X} ולכל $f \in \mathcal{H}$, אם מתקיימת הנחת ה-*realizability* ביחס ל- $\mathcal{H}, \mathcal{D}, f$ כאשר נריץ את A על $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ דוגמאות *i.i.d* מפילוג \mathcal{D} שתוייגו על ידי f אזי האלגוריתם החזיר היפותזה $h \in \mathcal{H}$ אשר:

$$\mathbb{P}_{S|x \sim \mathcal{D}^m} [L_{\mathcal{D},f}(h) \leq \epsilon] \geq 1 - \delta$$

או במילים אחרות ההסתברות שהאלגוריתם יחזיר היפותזה בעלת איבוד על כל \mathcal{D} שגבוה מהסף שקבענו (ϵ) קטן מ- δ . ϵ מכונה ה-*accuracy* ו- δ מכונה ה-*confidence*.

Sample Complexity 6.2

הפונקציה $m_{\mathcal{H}}$ מכונה ה-*sample complexity* של מחלקת ההיפותזות והיא מגלמת את הקשר בין ה-*accuracy* וה-*confidence* שברצוננו להשיג למספר הדוגמאות (*samples*) עליהם צריך להתאמן האלגוריתם. לרוב נרצה לקבוע ϵ, δ בהתאם לדרישות עבור פתרון הבעיה הלימודית ופונקציה זו למעשה נותנת חסם תחתון למספר הדוגמאות עליהם נתאמן.

טענה 6.2 יהיו $\epsilon, \delta \in (0, 1)$. אם $m \geq \left\lceil \frac{\log(|\mathcal{H}|/\delta)}{\epsilon} \right\rceil$ אזי לכל \mathcal{D}, f בהסתברות לפחות $1 - \delta$ מתקיים

$$L_{\mathcal{D},f}(ERM_{\mathcal{H}}(S)) \leq \epsilon$$

מסקנה 6.3 עבור \mathcal{H} מחלקת היפותזות סופית אזי \mathcal{H} הינה *PAC learnable* עם $m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log(|\mathcal{H}|/\delta)}{\epsilon} \right\rceil$.

הוכחה לטענה ניתן למצוא בספר *Understanding Machine Learning* בעמודים 39 – 40.

6.3 Agnostic PAC

כפי שראינו עד כה ב-*PAC learning* אנו מניחים את קיומה של פונקציית תיוג נכונה f שבמחלקת ההיפותזות. הנחה זו מגבילה אותנו שכן יתכן כי לא קיימת פונקציה כזו כלל או שהיא איננה במחלקת ההיפותזות. על כן נרצה להרחיב את המודל להתמודד עם מצב זה.

6.3.1 Data Generating Distribution

הנחת ה-*realizability* למעשה קבעה כי f זו מתייגת נכון בהסתברות 1 כל $\mathcal{X} \sim \mathcal{D}$. על מנת להשתחרר מהגבלה זו נגדיר את הפילוג על מרחב התיוגים בנוסף למרחב הדוגמאות: $\mathcal{X} \times \mathcal{Y} = \mathcal{Z} \sim \mathcal{D}$. נכנה את \mathcal{D} בתוך ה-*joint distribution* מעל נקודות ב-*domain* ותיוגים כאשר:

• $\mathcal{D}_{\mathcal{X}}$ הינה ההתפלגות מעל הנקודות ב-*domain* ומכונה *marginal distribution*.

• $\mathcal{D}((x, y) | x)$ הינה ההסתברות המותנית לתיוגים עבור נקודות ב-*domain*.

נשים לב כי אילו מחלקת היפותזות \mathcal{H} הינה *Agnostic PAC learnable* (אשר נראה הגדרה מלאה בהמשך) ונניח כעת כי אכן מתקיימת הנחת ה-*realizability* אזי למעשה הפונקציה המתייגת נכונה בהסתברות 1 תספק את ה- $y \in \mathcal{Y}$ הנכון ואת יתר התיוגים בהסתברות 0. על כן נוכל להסתכל על *PAC* כמקרה פרטי של *Agnostic PAC*.

6.3.2 True error/Empirical Risk

עבור $\mathcal{X} \times \mathcal{Y} = \mathcal{Z}$ אזי פונקציית האיבוד מוגדרת בתוך $l : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$ ועל כן ה-*risk function* הינו האיבוד המצופה ממסווג $h \in \mathcal{H}$ ביחס לפילוג \mathcal{D} מעל \mathcal{Z} :

$$L_{\mathcal{D}}(h) = \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} [l(h, \mathbf{z})]$$

ונגדיר את ה-*empirical risk* מעל \mathcal{Z}^m $S = (z_1, \dots, z_m) \subseteq \mathcal{Z}^m$ בתור:

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z_i)$$

הגדרת ה-*ERM*, שכן הינה על סט האימון, נשארת זהה:

$$ERM_{\mathcal{H}}(S) \in \operatorname{argmin}_{h \in \mathcal{H}} L_S(h) = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{m} |\{x_i \in S \mid h(x_i) \neq y_i\}|$$

כעת נראה את ההגדרה הפורמלית של *Agnostic PAC*:

הגדרה 6.4 נאמר כי מחלקת היפותזות \mathcal{H} הינה *Agnostic PAC learnable* אם קיימים פונקציה $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ ואלגוריתם למידה \mathcal{A} עם התכונה: שלכל $\epsilon, \delta \in (0, 1)$ ולכל התפלגות \mathcal{D} מעל $\mathcal{X} \times \mathcal{Y}$ כאשר מריצים את \mathcal{A} על $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ דוגמאות $i.i.d$ מפילוג \mathcal{D} אזי האלגוריתם מחזיר היפותזה h כך ש:

$$\mathcal{D}^m(\{S \in \mathcal{Z}^m : L_{\mathcal{D}}(h) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon\}) \geq 1 - \delta$$

6.4 דוגמה - Coin Prediction

כבעיית חימום לסיווג נתעסק בשערוך תוצאת הטלת מטבע. בהינתן משתנה מקרי z המייצג הטלת מטבע עם סיכוי הצלחה p נרצה ללמוד מהו p . המטבע הינו $z \sim \text{Ber}(p)$.

יהיו $z_1, \dots, z_m \sim \mathcal{D}_z$ מ"מ ש"ה ב"ת (כלומר $i.i.d$) סט אימון S אשר ההתפלגות שלו \mathcal{D}_z^m . אם כך הקלט לבעיית הלמידה \mathcal{A} הינו S (אקראי שנבחר בהתאם ל- \mathcal{D}^m והפלט, הפרדיקציה, הינה $\mathcal{A}(S)$ ונסמן \hat{p} .

היות ואופן התפלגות המידע אינו ידוע לנו נרה להכניס פרמטר *accuracy* למערכת ($\epsilon > 0$) ונרשה ל- \hat{p} לסטות מ- p בכלל היותר ϵ . בנוסף, יתכן כי סט האימון איננו מייצג את ההתפלגות הכללית ולכן נכניס *confidence* למערכת ($\delta > 0$) אשר יבטיח כי

$$\mathbb{P}(\{S : |\hat{p} - p| > \epsilon\}) < \delta$$

על כן באופן פורמלי אלגוריתם למידה \mathcal{A} עבור *coin prediction* הינו $f : \cup_{m=1}^{\infty} \{0, 1\}^m \rightarrow [0, 1]$ המקיימת את ההבאים: קיימת פונקציה $m : (0, 1)^2 \rightarrow \mathbb{N}$ כל שלכל $\epsilon, \delta \in (0, 1)$ ולכל התפלגות \mathcal{D} המתאימה ל- $p \in [0, 1]$, כאשר נריץ את \mathcal{A} על $m \geq m_{\mathcal{A}}(\epsilon, \delta)$ דוגמאות $i.i.d$ על פי \mathcal{D}^m האלגוריתם יחזיר \hat{p} כך שבהסתברות לפחות $1 - \delta$ נקבל כי $|\hat{p} - p| \leq \epsilon$.

בעיה מטילים מטבע בעל הסתברות p לקבלת H m פעמים. רוצים לשערך את p .

פתרון מרחב ההסתברות הינו $\Omega = \{H, T\}^m$ ו- $D^m = \frac{1}{2^m}$ (כלומר יוניפורמי). נגדיר z_1, \dots, z_m כאשר $z_i = \mathbb{1}_{w_i=H}$. האלגוריתם יקבל כקלט סט ערכים של m הטלות והפלט הינו $\bar{z} = \frac{1}{m} \sum_{i=1}^m z_i$. נשים לב כי $\mathbb{E}[\bar{z}] = p$ ומאי שיוויון הופדינג נקבל כי

$$\mathcal{D}^m[|\bar{z} - p| \geq \epsilon] \leq 2 \exp(-2m\epsilon^2)$$

נרצה כי ביטוי זה יהיה קטן מ- δ ולכן:

$$2 \exp(-2m\epsilon^2) \leq \delta \Leftrightarrow -2m\epsilon^2 \leq \ln\left(\frac{\delta}{2}\right) \Leftrightarrow \frac{1}{2\epsilon^2} \ln\left(\frac{2}{\delta}\right) \leq m$$

7 $VC - Dimension$

על אף שראינו כי סופיות של מחלקת היפותזות \mathcal{H} הינו תנאי מספיק ללמידה הוא איננו תנאי הכרחי. כדי לספק תנאי רחב יותר גם עבור מחלקות אינסופיות נזכר ברעיון ההוכחה של *No Free Lunch*. במקרה המשפט הראינו שללא הגבלת מחלקת ההיפותזות לכל אלגוריתם, באמצעות שיטת היריב נוכל ליצור התפלגות עליה אלגוריתם למידה אחד ייכשל ואחר שיצליח. בשביל לעשות זאת השתמשנו בתת קבוצה סופית של הדומיין $C \subset \mathcal{X}$. שיטת היריב החזירה אלגוריתם שעבור קבוצה סופית זו בחרה בפונקציית תיוג שתיכשל על קלט זה.

הגדרה 7.1 תהי \mathcal{H} מחלקת היפותזות מ- \mathcal{X} ל- $\{0, 1\}$ ותהי $C = \{c_1, \dots, c_m\} \subset \mathcal{X}$. נאמר כי ה-*restriction* של \mathcal{H} על C הינה תת קבוצת הפונקציות מ- C ל- $\{0, 1\}$ שניתן לקבל מ- \mathcal{H} :

$$\mathcal{H}_C = \{h(c_1), \dots, h(c_m) : h \in \mathcal{H}\}$$

הגדרה 7.2 נאמר כי מחלקת היפותזות \mathcal{H} מנתצת (*shatters*) קבוצה וסופית $C \subset \mathcal{X}$ אם ה-*restriction* של \mathcal{H} על C הינו סט כל הפונקציות מ- C ל- $\{0, 1\}$. כלומר $|\mathcal{H}_C| = 2^{|C|}$.

הגדרה 7.3 מימד ה- VC ($VC - Dimension$) של מחלקת היפותזות \mathcal{H} , מסומן ב- $VCdim(\mathcal{H})$, הינו הגודל המקסימלי של $C \subset \mathcal{X}$ אשר \mathcal{H} מנתצת את C .

משפט 7.4 תהי מחלקת היפותזות \mathcal{H} בעלת מימד VC אינסופי אזי \mathcal{H} איננה *PAC learnable*.

לצורך הוכחת המשפט קיימת טענת עזר (6.4 בספר *understanding machine learning*, עמוד 69) שאם \mathcal{H} מנתצת סט מגודל $2m$ לא נוכל ללמוד את \mathcal{H} באמצעות m דוגמאות בלבד.

Uniform Convergence 8

נראה תנאי, *uniform convergence*, אשר הינו תנאי מספיק ללמידה. ראינו כי *ERM* מחזיר היפותזה הממזערת שגיאה על האימון ובתקווה שכך גם על \mathcal{D} הכללי. על כן מספיק לבקש שה-*empirical risk* של כל $h \in \mathcal{H}$ מקרב מספיק ל-*true error*.

הגדרה 8.1 סט אימון S הינו ϵ -מייצג (*ϵ -representative*) (ביחס לדומיין \mathcal{Z} , מחלקת היפותזות \mathcal{H} , פונקציית איבוד l והתפלגות \mathcal{D}) אם

$$\forall h \in \mathcal{H} \quad |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon$$

למה 8.2 יהי סט אימון S אשר הינו $\frac{\epsilon}{2}$ -מייצג אזי לכל פלט של $ERM_{\mathcal{H}}(S)$ מתקיים כי:

$$h_S \in \operatorname{argmin}_{h \in \mathcal{H}} L_S(h) \quad L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon$$

הוכחה: תהי $h \in \mathcal{H}$ אזי:

$$L_{\mathcal{D}}(h_S) \stackrel{\frac{\epsilon}{2}-rep.}{\leq} L_S(h_S) + \frac{\epsilon}{2} \stackrel{min\ of\ group}{\leq} L_S(h) + \frac{\epsilon}{2} \stackrel{\frac{\epsilon}{2}-rep.}{\leq} L_{\mathcal{D}}(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2} = L_{\mathcal{D}}(h) + \epsilon$$

■

על כן נסיק כי עבור סט $\frac{\epsilon}{2}$ -מייצג עיקרון ה-*ERM* יחזיר היפותזה טובה.

הגדרה 8.3 מחלקת היפותזות \mathcal{H} הינה בעלת התכונה *uniform convergence* אם קיימת פונקציה $m_{\mathcal{H}}^{UC} : (0, 1)^2 \rightarrow \mathbb{N}$ כך שלכל $\epsilon, \delta \in (0, 1)$ ולכל התפלגות \mathcal{D} מעל \mathcal{Z} אז

$$\mathcal{D}^m(\{S \in \mathcal{Z}^m : S \text{ is } \epsilon\text{-representative}\}) \geq 1 - \delta$$

מסקנה 8.4 אם \mathcal{H} הינה *uniform convergence* עם $m_{\mathcal{H}}^{UC}$ אזי \mathcal{H} הינה *agnostic PAC learnable* עם *sample complexity* של $m_{\mathcal{H}} \leq m_{\mathcal{H}}^{UC}(\frac{\epsilon}{2}, \delta)$ ופרדיגמת *ERM* הינה *agnostic PAC learner* טובה עבור \mathcal{H} .

טענה 8.5 תהי \mathcal{H} מחלקת היפותזות סופית ופונקציית איבוד l חסומה ב- $[0, 1]$ אזי \mathcal{H} הינה *agnostic PAC learnable* באמצעות *ERM* ו-*sample complexity*:

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{2 \log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{\epsilon^2} \right\rceil$$

9 המשפט היסודי של הלמידה הסטטיסטית

ישנן שתי גרסאות שונות למשפט. הראשונה מספקת חסמים עבור פונקציית ה- $sample complexity$ והשנייה מספקת תנאים שקולים ללמידה.

משפט 9.1 \mathcal{H} תהי מחלקת היפותזות מדומיין \mathcal{X} ותיוגים $\{0, 1\}$ ותהי L^{0-1} פונקציית האיבוד. נניח כי $VCdim(\mathcal{H}) = d < \infty$ אזי קיימים קבועים אוניברסליים C_1, C_2 כך שמתקיים:

1. \mathcal{H} הינה PAC learnable עם

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

2. \mathcal{H} הינה $agnostic PAC$ learnable עם

$$C_1 \cdot \frac{d + \log\left(\frac{1}{\delta}\right)}{\epsilon^2} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \cdot \frac{d + \log\left(\frac{1}{\delta}\right)}{\epsilon^2}$$

3. \mathcal{H} בעלת תכונת ה- $uniform convergence$ עם

$$C_1 \cdot \frac{d + \log\left(\frac{1}{\delta}\right)}{\epsilon^2} \leq m_{\mathcal{H}}^{UC}(\epsilon, \delta) \leq C_2 \cdot \frac{d + \log\left(\frac{1}{\delta}\right)}{\epsilon^2}$$

(*) רק המקרה הראשון הוצג בכיתה.

משפט 9.2 \mathcal{H} תהי מחלקת היפותזות מדומיין \mathcal{X} ותיוגים $\{0, 1\}$ ותהי L^{0-1} פונקציית האיבוד. התנאים הבאים שקולים:

1. \mathcal{H} הינה $uniform convergence$.

2. כל היפותזות ERM הינה $agnostic PAC$ learner טובה עבור \mathcal{H} .

3. \mathcal{H} הינה $agnostic PAC$ learnable.

4. \mathcal{H} הינה PAC learnable.

5. כל היפותזות ERM הינה PAC learner טובה עבור \mathcal{H} .

6. $VCdim(\mathcal{H}) < \infty$.

10 הבנת Bias – Complexity Tradeoff

נראה כי כאשר בוחרים מחלקת היפותזות לבעיית למידה נרצה למעשה לנסות להשיג 2 דברים:

1. נרצה שמחלקה זו תכיל פונקציית תיוג נכונה או לכל הפחות שהפונקציה בעלת האיבוד הנמוך ביותר במחלקה אכן נמוך מספיק לצורכינו. על כן נרצה להגדיל את המחלקה.
2. לא נוכל לבחור מחלקה רחבה מדי שכן אחרת מ-"*No Free Lunch*" נקבל כי לא נצליח ללמוד את הדומיין. על כן נרצה להקטין את המחלקה.

10.1 Error Decomposition

כדי לראות כיצד להשיג 2 מטרות אלה נפרק את השגיאה של לומד $ERM_{\mathcal{H}}$ לשני מרכיבים. עבור h_S היפותזת $ERM_{\mathcal{H}}$ נאמר:

$$L_{\mathcal{D}}(h_S) = \underbrace{\epsilon_{app}}_{bias} + \underbrace{\epsilon_{est}}_{variance as in complexity}$$

$$\epsilon_{app} = \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$$

$$\epsilon_{est} = L_{\mathcal{D}}(h_S) - \epsilon_{app}$$

כאשר:

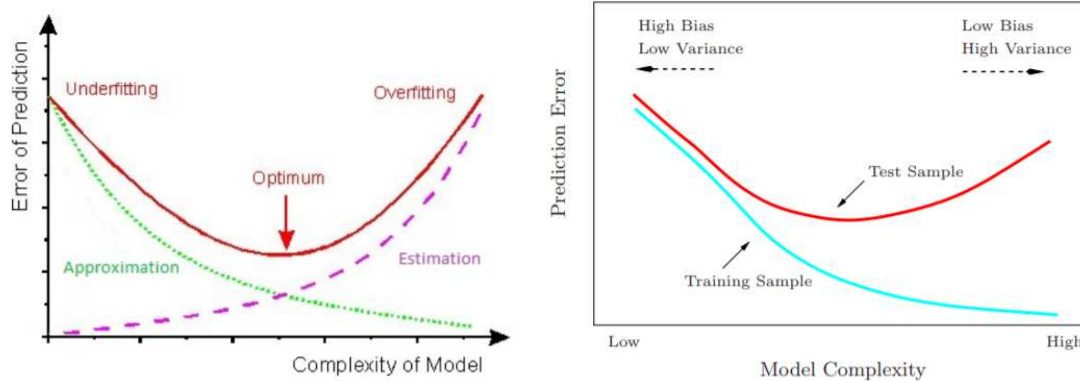
1. The Approximation Error - השגיאה המינימלית שיכולה להתקבל על ידי פונקציה במחלקה. למעשה שגיאה זו מייצגת את כמות ה"סיכון" בהגבלתנו את מחלקת ההיפותזות (כמה עשירה המחלקה). כלומר ה-"*inductive bias*" שקיים במערכת. שגיאה זו למעשה איננה תלויה בגודל סט האימון אלא רק במחלקת ההיפותזות. במקרה שמתקיימת הנחת ה"*realizability*" שגיאה זו הינה 0. באופן כללי *bias* מתייחס להעדפה של דבר אחד על פני אחר כתוצאה מדעה מוקדמת.
2. The Estimation Error - המרחק בין שגיאת ה-"*approximation*" לבין השגיאה שמתקבלת על ידי הפרדיקציה של $ERM_{\mathcal{H}}$. שגיאה זו נובעת מכך שה-"*empirical risk*" הינה רק שערור של השגיאה האמיתית על \mathcal{D} . איכות שגיאה זו תלויה בגודל סט האימון ובגודל (או $VCdim$) מחלקת ההיפותזות. נוכל להתייחס לגודל של \mathcal{H} כמדד ל-"*complexity*" שלה.

היות ובאופן כללי נרצה לצמצם את השגיאה הכוללת אנו ניצבים בפני *tradeoff* בין שני אלה. עבור מחלקה עשירה מאוד נוריד את שגיאת ה-"*approximation*" אבל נסתכן בהגדלת שגיאת ה-"*estimation*" ונסתכן ב-"*overfitting*" (הליכה אחרי הרעש האינרנטי בדגימות). עבור מחלקה מצומצמת נוריד את שגיאת ה-"*estimation*" אבל נסתכן בהגדלת שגיאת ה-"*approximation*" וב-"*underfitting*".

10.2 Tradeoff

ראינו כי ה-*bais* מודד את המרחק בין החיזוי ה"ממוצע" לאמת וה-*variance* מודד כמה וריאבילית הפרדיקציה.

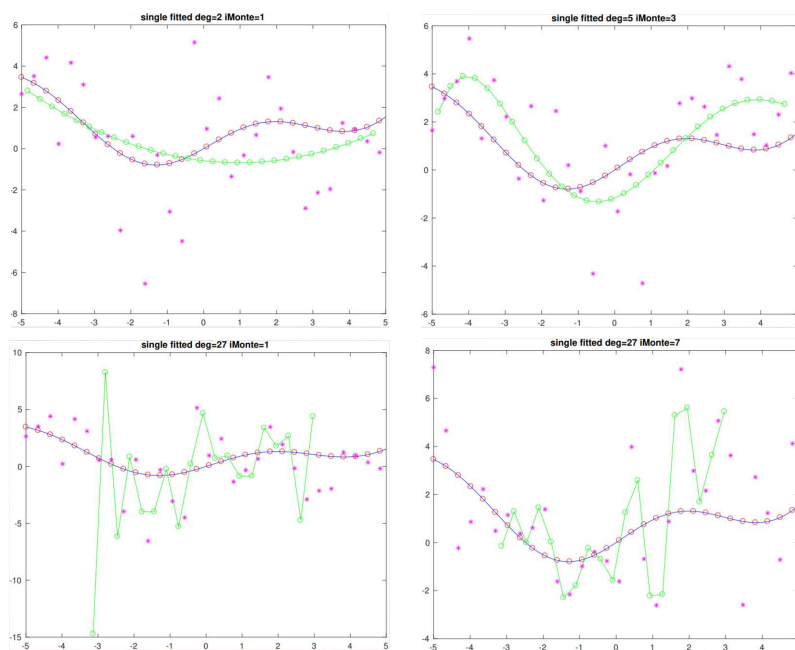
- יותר מידע \Leftarrow פחות *variance*.
- מודל יותר *complicated* \Leftarrow יותר *variance*.
- יותר רעש σ^2 \Leftarrow יותר *variance*.



ברמת האינטואיציה עבור מודל פשוט יהיה לנו *bias* גבוה ו-*variance* נמוך מכיוון מוחלקת ההיפותזות פשוטה ונקבל שגיאה גבוהה. עבור מודל מורכב נקבל *bias* נמוך ו-*variance* גבוהה מכיוון ו"נלך אחרי הרעש" שבדגימות ונקבל שגיאה גבוהה. על כן לכל מודל נצטרך להבין מהם המנופים שנוכל לשחק בהם כדי למצוא את נקודת המינימום שבסט ה-*test*.

10.3 דוגמה עבור פולינומים

בתמונות הבאות: הנקודות הירוקות - החיזוי, הנקודות אדומות - ערך אמת, כוכביות - דוגמאות האימון כאשר ההבדל בין לבין האדומות נובע מהכנסת רכיב הרעש, הקו הירוק - המשוואה, הקו הכחול המשוואה ה"מושלם".



מה שאנו רואים זה שפולינומים ממעלה נמוכה לא מספקים פרדיקציה טובה מספיק (*underfitting*) בעוד שבמעלות גבוהות המסווג "עונה כל כך טוב" על הדוגמאות שלמעשה הוא מסווג את הרעש ולא את ה-*data* שמאחורי הרעש. המסווג מאוד "גמיש" והוא "רץ אחרי הרעש". מקרה זה מתאר *overfitting*.

חלק IV

תיאוריה לאלגוריתמים

בחלק זה נראה אוסף אלגוריתמים הפותרים בעיות למידה. חלק מהאלגוריתמים מספקים פתרון לעולם תוכן מסויים (לדוגמה רגרסיה לינארית) בעוד שאחרים הינם אלגוריתמים לשיפור תכונה כללית מסוימת (כמו *boosting*).

11 מסווגים לינאריים

הגדרה 11.1 מחלקת הפונקציות האפיניות הינה $L_d = \{h_{w,b} : w \in \mathbb{R}^d, b \in \mathbb{R}\}$ כאשר $h_{w,b}(x) = \langle w, x \rangle + b = \left(\sum_{i=1}^d w_i x_i\right) + b$. נוח לסמן בתור $L_d = \{x \mapsto \langle w, x \rangle + b : w \in \mathbb{R}^d, b \in \mathbb{R}\}$.

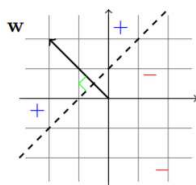
במקרים רבים נוהגים לשלב את b , הנקרא גם ה-*bias*, לתוך w בקואורדינטה ה-0 ולהוסיף לכל $x \in \mathcal{X}$ קואורדינטה ראשונה עם הערך 1. כלומר $x' = (1, x_1, \dots, x_d) \in \mathbb{R}^{d+1}$, $w' = (b, w_1, \dots, w_d) \in \mathbb{R}^{d+1}$ ומכאן כי: $h_{w,b}(x) = \langle w, x \rangle + b = \langle w', x' \rangle$.

11.1 חצאי מרחב (*Halfspaces*)

נתבונן במחלקת ההיפותוזות של מפרידים לינאריים עבור בעיות קלסיפיקציה $|\mathcal{Y}| = 2$ לרוב $\mathcal{Y} = \{\pm 1\}$ או $\mathcal{Y} = \{0, 1\}$:

$$HS_d = \text{sign} \circ L_d = \{x \mapsto \text{sign}(h_{w,b}(x)) : h_{w,b} \in L_d\}$$

במקרה של $d = 2$ היפותזה המוחזרת מהמחלקה מגדירה על-מישור מפריד אשר מאונך ל- w . דגימות (ערכי קלט x) מעל העל-מישור יקבלו +, ודגימות מתחת יקבלו -.



11.1.1 פתרון ERM עבור *halfspaces* באמצעות *Perceptron*

אלגוריתם ה-*Perceptron* עובד באופן איטרטיבי ליצירת סדרת וקטורים $w^{(1)}, w^{(2)}, \dots$ כאשר בכל איטרציה t האלגוריתם מוצא *sample* x_i אשר $w^{(t)}$ טעה $(\text{sign}(\langle w, x_i \rangle) \neq y_i)$ ומעדכן את $w^{(t)}$ על ידי הוספת $x_i y_i$. אנו מקבלים:

$$y_i \langle w^{(t+1)}, x_i \rangle = y_i \langle w^{(t)} + y_i x_i, x_i \rangle = y_i \langle w^{(t)}, x_i \rangle + \|x_i\|^2$$

כלומר כאשר מעדכנים את הוקטור w מכוונים אותו מעט לכיוון ה-*sample* i .

Batch Perceptron

input: A training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

initialize: $\mathbf{w}^{(1)} = (0, \dots, 0)$

for $t = 1, 2, \dots$

if $(\exists i \text{ s.t. } y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle \leq 0)$ **then**

$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + y_i \mathbf{x}_i$

else

output $\mathbf{w}^{(t)}$

משפט 11.2 יהיו $S = \{(x_i, y_i)\}_{i=1}^m$ ברי הפרדה לינארית, $R = \max_i \|x_i\|$ ו- $B = \min \{\|w\| : \forall i \in [m] y_i \langle w, x_i \rangle \geq 1\}$ אזי האלגוריתם יעצור אחרי לכל היותר $(RB)^2$ איטרציות ובסופן $\forall i \in [m] y_i \langle w^{(t)}, x_i \rangle > 0$.

הערה 11.3 אלגוריתם *Perceptron* בהכרח יתכנס, אולם קצב ההתכנסות תלוי ב- B שבמקרים מסוימים עשוי להיות אקספוננציאלי ב- d .

11.1.2 פתרון ERM עבור *halfspaces* באמצעות תכנון לינארי

עבור המקרה של חצאי-מרחב הומוגניים (נוכל להרחיב ללא הומוגניים באופן דומה למתואר בתחילת הפרק) בהינתן סט אימון S אזי היפותזת ERM הינה וקטור $w \in \mathbb{R}^d$ אשר:

$$\forall i \in [m] \text{ sign}(\langle w, x_i \rangle) = y_i \iff y_i \cdot \langle w, x_i \rangle > 0$$

נציג זאת בצורה סטנדרטית של תכנון לינארי. תחת הנחת הריאלזביליות קיים w^* המקיים את כל התנאים אזי נוכל להניח כי

$$\forall i \in [m] y_i \langle w, x_i \rangle \geq 1$$

אזי נגדיר

$$\gamma = \min_i \{y_i \langle w, x_i \rangle\}$$

ועבור $\bar{w} = \frac{1}{\gamma} w^*$ מתקיים:

$$\forall i \in [m] y_i \langle \bar{w}, x_i \rangle = \frac{y_i \langle w^*, x_i \rangle}{\gamma} \geq 1$$

ובצורה מטריציונית נקבל:

$$y_i \langle w, x_i \rangle = y_i x_i^T w$$

$$\begin{bmatrix} - & y_1 x_1^T & - \\ & & \\ - & y_m x_m^T & - \end{bmatrix} w \geq 1_m$$

11.1.3 $VCdim$ של $halfspaces$

משפט 11.4 ה- $VCdim$ של מחלקת ההיפותזות של ה- $halfspaces$ ההומוגניים מעל \mathbb{R}^d הוא d .

הוכחה:

1. תהי $\{e_1, \dots, e_d\} = C \subset \mathcal{X}$ קבוצת וקטורי הבסיס הסטנדרטי $((e_i)_j = \mathbb{1}_{i=j})$. נראה כי לכל תיוג $y_1, \dots, y_d \in \{-1, 1\}^d$ נבחר $w = (y_1, \dots, y_d)^T$ ואז $\langle w, e_i \rangle = y_i$ ו- $\forall i \in [d]$ על כן \mathcal{H} מנתצת את C ומתקיים $VCdim(\mathcal{H}) \geq d$.

2. תהי $\{x_1, \dots, x_{d+1}\} = C \subset \mathcal{X}$. היות ו- $dim(\mathbb{R}^d) = d$ קבוצה זו תלויה לינארית ולכן קיימים $a_1, \dots, a_{d+1} \in \mathbb{R}$ כולם אפס כך ש: $\sum a_i x_i = 0$. נסמן $J = \{i : a_i < 0\}$, $I = \{i : a_i > 0\}$. בהכרח אחת הקבוצות לא ריקה.

• אם שתי הקבוצות לא ריקות אזי מתקיים $\sum_{i \in I} a_i x_i = \sum_{j \in J} |a_j| x_j$. נניח בשלילה כי \mathcal{H} מנתצת את C ($|C| = d+1$). אזי קיים w כך ש: $\forall i \in I \langle w, x_i \rangle > 0 \wedge \forall j \in J \langle w, x_j \rangle < 0$. מכאן נובע כי:

$$0 < \sum_{i \in I} a_i \langle w, x_i \rangle = \left\langle w, \sum_{i \in I} a_i x_i \right\rangle = \left\langle w, \sum_{j \in J} |a_j| x_j \right\rangle = \sum_{j \in J} |a_j| \langle w, x_j \rangle < 0$$

בסתירה.

• אם אחת הקבוצות ריקה, נניח בה"כ J , נקבל בהתאם:

$$0 < \sum_{i \in I} a_i \langle w, x_i \rangle = \left\langle w, \sum_{i \in I} a_i x_i \right\rangle = \left\langle w, \sum_{j \in J} |a_j| x_j \right\rangle = \langle w, 0 \rangle = 0$$

בסתירה.

משפט 11.5 חה- $VCdim$ של מחלקת ההיפותזות של ה- $halfspaces$ הלא הומוגניים מעל \mathbb{R}^d הוא $d+1$.

הוכחה: תחילה נוודא ניתוח של הקבוצה $\{0, e_1, \dots, e_d\} = C \subset \mathcal{X}$ כשם שהוכחנו מקודם. לאחר מכן, באופן זהה להוכחה לעיל, נראה שלא ניתן לנתח קבוצה $\{x_1, \dots, x_{d+2}\} = C \subset \mathcal{X}$. נסיק כי $VCdim(HS_d) = d+1$ כנדרש.

11.2 רגרסיה לינארית

11.2.1 פתרון הבעיה

נציג תחילה עולם בעיה שנרצה ללמוד ונראה כיצד בעזרת המונחים השונים נתן אלגוריתם למידה עבור הבעיה. נניח שאנו סוחרי נדלן וברצוננו לשערך באיזה מחיר נוכל למכור נכס מסויים. עבור נכס ספציפי נוכל להתייחס לגודלו, למספר החדרים, למספר חדרי השינה, האם בבניין או בית פרטי, באיזה קומה, האם יש נוף, מצב הנכס ופרמטרים רבים נוספים. כל אלה הם ה- $features$ שבעזרתם ננסה לשערך את המחיר $label$.

על כן בעולם בעיה זו ה- $domain$ שלנו הוא נכסים, ה- $features$ הם התכונות שתיארנו וה- $labels$ הם מחיר הנכס. נשים לב שהתיוגים חיים ב- \mathbb{R} - ולכן מדובר פה בבעיית למידה של רגרסיה. באופן פורמלי: $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \mathbb{R}$ ואנו מחפשים $f: \mathcal{X} \rightarrow \mathcal{Y}$.

כעת נרצה לתאר את מחלקת ההיפותזות. בהתאם לעיקרון ה- $no free lunch$ עלינו להניח דבר מה על מחלקת ההיפותזות. בצורת למידה זו אנו מניחים כי התיוגים \mathcal{Y} , בדוגמה הנ"ל מחירי הנכסים, הינם לינאריים ב- $features$ אזי החלקת ההיפותזות הינה מהצורה:

$$\mathcal{H}_{reg} = \left\{ (x_1, \dots, x_d) \mapsto w_0 + \sum_{i=1}^d x_i w_i \mid w_0, w_1, \dots, w_d \in \mathbb{R} \right\}$$

בנוסף נניח כי מתקיימת $the realizability assumption$ אזי קיימת $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$ אשר מתייגת נכונה את הדגימות ונניח כי $f \in \mathcal{H}_{reg}$. בהגדרה זו כל $feature$ מקבל משקל w_i וקיים w_0 אשר הוא ה- $intercept$ - ההסטה של הישרייה מראשית הצירים. למידת בעיה זו משמע מציאת הוקטור $w \in \mathbb{R}^{d+1}$.

כעת נבין כיצד נראה ה- $data$ עליו אנו עובדים. תחילה לכל דגימה ב- \mathcal{X} נראה כי היא למעשה וקטור $(x_1, \dots, x_d) \in \mathbb{R}^d$. היות וקיים ה- $intercept$ נבחר לתאר את הדגימות במימד גבוה יותר: $(1, x_1, \dots, x_d) \in \mathbb{R}^{d+1}$. לכן עבור m דגימות ($samples$) ה- d $features$ נוכל לארגן את סט האימון במטריצה:

$$X = \begin{bmatrix} | & & | \\ x_1 & \dots & x_m \\ | & & | \end{bmatrix}, x_i \in \mathbb{R}^{d+1}$$

הגדרה 11.6 תהי $X \in \mathbb{R}^{d \times m}$ ופירוק SVD $X = U \Sigma V^T$ אזי ה- $Moore-Penrose pseudo inverse$ של x הוא:

$$X^\dagger = V \Sigma^\dagger U^T \quad \Sigma_{i,i}^\dagger = \begin{cases} (\Sigma_{i,i})^{-1} & \Sigma_{i,i} \neq 0 \\ 0 & \Sigma_{i,i} = 0 \end{cases}$$

טענה 11.7 לימוד בעיית הרגרסיה הלינארית שקול למציאת פתרון w למערכת המשוואות $XX^T w = Xy \iff X^T w = y$.

הוכחה: נראה כי בהתאם למטריצה XX^T נוכל לחלק למקרים עבורם נמצא את הפתרון למערכת המשוואות: תחילה נניח כי XX^T הפיכה ונוכיח כי קיים למערכת המשוואות פתרון יחיד וחילוצו: $\hat{w} = (XX^T)^{-1} Xy = (X^\dagger)^T y$. היותר ו- XX^T הפיכה קיימת לה מטריצה הפוכה. על כן $XX^T w = Xy \iff w = (XX^T)^{-1} Xy$ יהי פירוק SVD .

של $X: X = U\Sigma V^T$. היות ו- XX^T הפיכה קיים לה פירוק EVD . נוכל לבחור כי המטריצה האורתונורמלית תהיה U מפירוק ה- SVD של X ולכן $(XX^T)^{-1} = UD^{-1}U^T$. מכאן:

$$\begin{aligned}(XX^T)^{-1}X &= (UD^{-1}U^T)(U\Sigma V^T) \\ &= UD^{-1}U^T U \Sigma V^T \\ &= UD^{-1}\Sigma V^T \\ &\stackrel{(*)}{=} U\Sigma^\dagger V^T \\ &= \left(V(\Sigma^\dagger)^T U^T\right)^T \\ &= (V\Sigma^\dagger U^T)^T \\ &= (X^\dagger)^T\end{aligned}$$

כאשר $(*)$ בגלל שלכל $i \neq j$ מתקיים כי $(D^{-1}\Sigma)_{i,j} = 0$ וכאשר $i = j$ אזי

$$(D^{-1}\Sigma)_{i,i} = \sum_k (D_{i,k}^{-1} \cdot \Sigma_{k,i}) = D_{i,i}^{-1} \Sigma_{i,i} = \frac{\sigma_i}{\sigma_i^2} = \sigma_i^{-1}$$

ולכן פתרון מערכת המשוואות הוא $\hat{w} = (X^\dagger)^T y$

כעת נניח כי XX^T איננה הפיכה ונוכיח כי קיימים אינסוף פתרונות והמינימלי והוא: $\|\hat{w}\| = \min\{\|w\| : X^T w = y\}$. היות XX^T איננה הפיכה אזי $\dim(Ker(XX^T)) \neq 0$ אזי מתקיים גם כי $\dim(X) \neq 0 \wedge \dim(X^T) \neq 0$ היות והגרעין איננו טריוויאלי למערכת המשוואות $X^T w = y$ קיימים אינסוף פתרון או אין פתרון כלל. יהי $y \perp Ker(X)$ אזי $y \in Im(X^T)$ ולכן $y \in Ker(X)^\perp$ מכאן כי קיים $w \in \mathbb{R}^d$ עבורו מתקיים $X^T w = y$ ונובע כי קיימים אינסוף פתרונות למערכת המשוואות.

על מנת לחלץ את הפתרון המינימלי יהיו $X = U\Sigma V^T$, $XX^T = UDU^T$ פירוקי SVD , EVD של המטריצות בהתאמה אזי:

$$\begin{aligned}X^T w &= y \\ XX^T w &= Xy \\ UDU^T w &= U\Sigma V^T y \\ \begin{bmatrix} d_1 & & & \\ & \ddots & & \\ & & d_r & \\ & & & 0 \end{bmatrix} U^T w &= \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & \\ & & & 0 \dots 0 \end{bmatrix} V^T y \\ D^\dagger D U^T w &= D^\dagger \Sigma V^T y\end{aligned}$$

כעת נשים לב כי:

$$D^\dagger D = \begin{bmatrix} d_1^{-1} & & & \\ & \ddots & & \\ & & d_r^{-1} & \\ & & & 0 \end{bmatrix} \begin{bmatrix} d_1 & & & \\ & \ddots & & \\ & & d_r & \\ & & & 0 \end{bmatrix} = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & 0 \end{bmatrix}$$

$$D^\dagger \Sigma = \begin{bmatrix} d_1^{-1} & & & \\ & \ddots & & \\ & & d_r^{-1} & \\ & & & 0 \end{bmatrix} \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & \\ & & & 0 \dots 0 \end{bmatrix} = \begin{bmatrix} \sigma_1^{-1} & & & \\ & \ddots & & \\ & & \sigma_r^{-1} & \\ & & & 0 \dots 0 \end{bmatrix} = (\Sigma^\dagger)^T$$

לכן נקבל כי:

$$U^T w = (\Sigma^\dagger)^T V^T y$$

$$\hat{w} = U (\Sigma^\dagger)^T V^T y$$

■

מסקנה 11.8 למשוואות הנורמליות $XX^T w = Xy$ בהכרח קיים פתרון (יחיד או אינסוף)

11.2.2 קיומו של רעש

כעת לאחר שראינו כיצד למצוא פתרון ללמידת רגרסיה לינארית נזכר כי בסט האימון שלנו קיים רעש ולכן למעשה הוא מהצורה $S = \{(x_i, f(x_i) + z_i)\}_{i=1}^m$ כאשר $z_1, \dots, z_m \stackrel{i.i.d}{\sim} (0, \sigma^2)$. כלומר מ"מ ב"ת ש"ה בעלי תוחלת 0 ושונות σ^2 . על כן כאשר נחפש את הפתרון נחפש $X^T w + z = y$ ומכאן כי $y \notin \text{Im}(X^T)$ ולמערכת ההומוגנית אין פתרון.

על כן בהתאם לעיקרון ERM נתבונן בפונקציית האיבוד $Least\ squares$. נראה כי מתקיים:

$$\sum_{(x_i, y_i) \in X, y} (y_i - X_i^T w)^2 = \|y - X^T w\|^2 = \langle y - X^T w, y - X^T w \rangle$$

התוצאה של $y_i - X_i^T w$ נקראת ה- i residual ומכאן במרה זה ה- ERM הינו $Residual Sum of Squares$:

$$\begin{aligned}
 RSS(w) &= \|y - X^T w\|^2 \\
 &\Downarrow \\
 \frac{\partial}{\partial w_i} RSS(w) &= -2 \sum_{i=1}^m (x_i) \cdot (y_i - X_i^T w) = 0 \\
 &\Downarrow \\
 \nabla RSS(w) &= -2X(y - X^T w) = 0 \\
 &\Downarrow \\
 X(y - X^T w) &= 0 \\
 X^T y &= X X^T w
 \end{aligned}$$

ומכאן כי אסטרטגיה זו אכן ממזערת את הטעות ומקיימת את עיקרון ה- ERM כנדרש.

Numerically Stable 11.2.3

נראה כי עולם הבעיה שתיארנו הינו ב- \mathbb{R}^d . היות וברצוננו לבטא כל זאת במחשב, ומחשב מסוגל לשמור רק מספר סופי (מצומצם) של ערכים אחרי הנקודה העשרונית (*floating point*) נרצה להתאים את המודל כך שיתמודד טוב יותר עם חוסר דיוק זה. לצורך ההבנה יתכן כי המחשב אינו יבדיל בין 0.00000 לבין 0.00001 ומכאן תנבא ההשפעה על תוצאת האלגוריתם. לכן עבור $\epsilon > 0$ שנבחר נגדיר את ה- *Moore-Penrose pseudo inverse* בתור:

$$\Sigma_{i,i}^{\dagger, \epsilon} = \begin{cases} (\Sigma_{i,i})^{-1} & \Sigma_{i,i} > \epsilon \\ 0 & \Sigma_{i,i} \leq \epsilon \end{cases}$$

Bias – Variance Decomposition 11.2.4

אחרי שאימנו מסווג \hat{w} יהיו $\tilde{x}_1, \dots, \tilde{x}_k$ דגימות חדשות אשר התיוגים האמתיים שלהן $\tilde{y}_1, \dots, \tilde{y}_k$. סיווג במסווג מתבצע על ידי $\hat{\tilde{y}} = \tilde{X}^T \hat{w}$. על מנת לבדוק את איכות הסיווג:

$$R(\tilde{y}, \hat{w}) = \mathbb{E} \left(\sum_{i=1}^k (\tilde{y}_i - \tilde{x}_i^T \hat{w})^2 \right) = \mathbb{E} \| \tilde{y} - \tilde{X}^T \hat{w} \|^2$$

כאשר מקור המקריות (ובגלל זה תוחלת) הוא הרעש שבתיגוים y - ומכאן גם הרעש ב \hat{w} . כעת כאשר נבדוק את ביטוי זה נקבל:

$$\begin{aligned}\mathbb{E} \left\| \tilde{y} - \tilde{X}^T \hat{w} \right\|^2 &= \mathbb{E} \left\| \tilde{y} - \hat{\tilde{y}} \right\|^2 \\ &= \mathbb{E} \left\| \tilde{y} - \mathbb{E} \tilde{y} + \mathbb{E} \tilde{y} - \hat{\tilde{y}} \right\|^2 \\ &= \mathbb{E} \left\| \tilde{y} - \mathbb{E} \tilde{y} \right\|^2 + \mathbb{E} \left\| \mathbb{E} \tilde{y} - \hat{\tilde{y}} \right\|^2 + 2\mathbb{E} \left\langle \tilde{y} - \mathbb{E} \tilde{y}, \mathbb{E} \tilde{y} - \hat{\tilde{y}} \right\rangle \\ &= \left\| \tilde{y} - \mathbb{E} \tilde{y} \right\|^2 + \sum_{i=1}^m \text{Var}(\hat{y}_i)\end{aligned}$$

אזי פירקנו את ה- $risk$ ל- $bias$ (משמאל) ול- $variance$ (מימין).

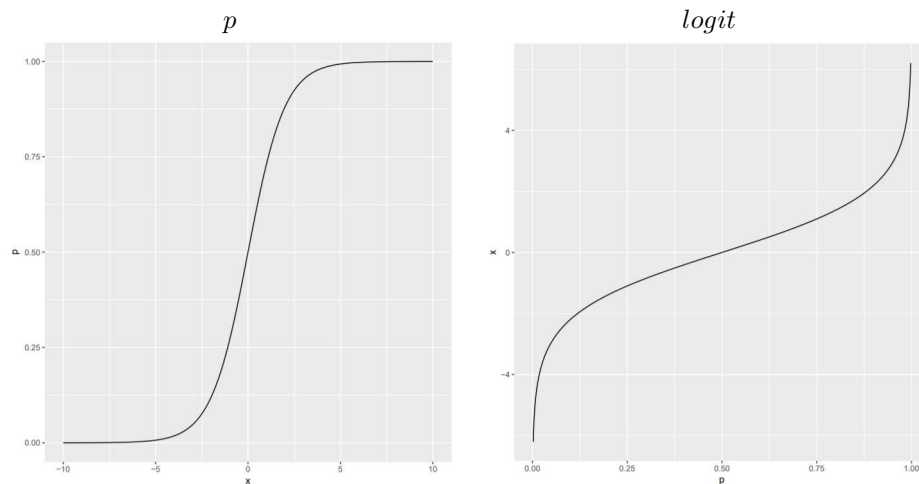
11.3 רגרסיה לוגיסטית

11.3.1 אינטואיציה

מודל הרגרסיה הלוגיסטית נובע מהרצון לשערך הסתברויות של סיווגים שונים (ב- $classification$ או $multiclassification$) באמצעות פונקציות לינאריות. במילים אחרות את ה- $likelihood$ של דגימה x להיות מסווגת k . במקרה של $classification$ אשר התיגוים הם $\mathcal{Y} = \{0, 1\}$ אזי נשאל $\mathbb{P}(Y_i | X_i = x_i)$ כאשר היות ומדובר בקלסיפיקציה אזי עבור $X_i = x_i$ נקבל כי $Y_i | X_i = x_i \sim \text{Ber}(p_i)$ ונרצה טרנספורמציה מתאימה מהפונקציה הלינארית ל- $[0, 1]$.

פונקציות קבות משמשות לרגרסיה הלוגיסטית כאשר מחלקה זו היא של פונקציות סיגמואידיות (מקבלות צורת S) החסומות ב- $[0, 1]$. אנו נתבונן בפונקציית ה- $\log - odds$ או ה- logit $\log - odds$ $\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$ כעת כשנסתכל על הפונקציה ההופכית שלה נקבל: $p(x) = \frac{e^x}{1+e^x}$ תחת ההנחה כי הפונקציה logit לינארית בתיגוים נקבל כי:

$$\begin{aligned}\text{logit}(y_i) &= ax_i + b \\ \Downarrow \\ y_i &= \frac{e^{ax_i+b}}{1 + e^{ax_i+b}}\end{aligned}$$



נראה כי כאשר ערכי x גדולים מאוד אזי p שואפת לתיוג 1 וכאשר ערכי x קטנים אזי p שואפת לתיוג 0.

11.3.2 הגדרה פורמלית

כעת במקום להניח כי $\mathbb{E}[Y_i | X_i = x_i]$ לינארית ב- x נניח כי לינארית ב- $g(x)$, כאשר g תהיה הפונקציה הלוגיסטית, ולכן:

$$\text{logit}(p_i) = \beta_0 + x_i^T \beta, \quad Y_i \sim \text{Ber}(p_i)$$

\Downarrow

$$\log \left(\frac{\mathbb{P}(Y_i = 1 | X_i = x_i)}{\mathbb{P}(Y_i = 0 | X_i = x_i)} \right) = \beta_0 + x_i^T \beta$$

כדי להגדיר את פונקציית ה- likelihood :

$$p_i = \frac{\exp(\beta_0 + x_i^T \beta)}{1 + \exp(\beta_0 + x_i^T \beta)}$$

ואז ה- likelihood ב- β_0, β על כל הדוגמאות היא:

$$L(\beta_0, \beta) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}$$

פונקציית ה- log-likelihood היא:

$$\ell(\beta_0, \beta) = \sum_{i=1}^n y_i (\beta_0 + x_i^T \beta) - \sum_{i=1}^n \log(1 + \exp(\beta_0 + x_i^T \beta))$$

וכעת על מנת לסווג נרצה למצוא את הסיווג שהכי likely . כלומר נחפש את הסיווג אשר $\arg\max_{\beta_0, \beta} \ell(\beta_0, \beta)$. בשונה מרגרסיה לינארית אין ביטוי סגור לחישוב זה אבל התוצאה נמצאת בתוך הקמור (convex) של β_0, β .

11.3.3 *multiclassification* מקרה

כאשר נרצה לתייג באוסף תיוגים למשל $\mathcal{Y} = \{1, \dots, K\}$ עבור $k \in \mathbb{N}$ אזי יהיו $\overline{\beta_1}, \dots, \overline{\beta_{K-1}}, \beta_1, \dots, \beta_{K-1} \in \mathbb{R}^d$

$$\begin{aligned} \log \frac{\mathbb{P}(G=1|X=x)}{\mathbb{P}(G=K|X=x)} &= \overline{\beta_1} + \beta_1^T x \\ &\vdots \\ \log \frac{\mathbb{P}(G=K-1|X=x)}{\mathbb{P}(G=K|X=x)} &= \overline{\beta_{K-1}} + \beta_{K-1}^T x \end{aligned}$$

12 *Bayes Classifiers*

12.1 *Bayes Optimal*

בהינתן התפלגות \mathcal{D} מעל $\mathcal{X} \times \{0, 1\}$ פונקציית התיוג האופטימלית, הנקראת גם *Bayes Optimal*, הינה

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \mathbb{P}[y=1|x] \geq \frac{1}{2} \\ 0 & \text{else} \end{cases}$$

טענה 12.1 מסווג ה-*optimal Bayes*, $f_{\mathcal{D}}$, הינו האופטימלי. כלומר $L_{\mathcal{D}}(f_{\mathcal{D}}) \leq L_{\mathcal{D}}(g)$ $\forall \mathcal{X} \rightarrow \{0, 1\}$.

הוכחה: יהי $x \in \mathcal{X}$ ותהי p ההסתברות של x להיות מתוייג $1 \in \mathcal{Y}$. לפי הגדרה מתקיים כי האיבוד של $h \in \mathcal{H}$:

$$L_{\mathcal{D}}(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathbb{1}_{h(x) \neq y}]$$

$$= \mathbb{E}_{x \sim \mathcal{D}} [\mathbb{E}_{y \sim D_{\mathcal{Y}|x}} [\mathbb{1}_{h(x) \neq y} | X = x]]$$

בפרט עבור $f_{\mathcal{D}}$ מתקיים כי ההסתברות של $\mathbb{1}_{f_{\mathcal{D}}(x) \neq y} | X = x$ היא:

$$\mathbb{P}(\mathbb{1}_{f_{\mathcal{D}}(x) \neq y} | X = x) = \mathbb{1}_{p \geq \frac{1}{2}} \cdot \mathbb{P}(Y = 1 | X = x) + \mathbb{1}_{p < \frac{1}{2}} \cdot \mathbb{P}(Y = 0 | X = x)$$

$$= \mathbb{1}_{p \geq \frac{1}{2}} \cdot (1 - p) + \mathbb{1}_{p < \frac{1}{2}} \cdot p$$

$$= \min(p, 1 - p)$$

מצד שני לכל $h \in \mathcal{H}$ ההסתברות של $\mathbb{1}_{f_D(x) \neq y} | X = x$ היא:

$$\begin{aligned} \mathbb{P}(\mathbb{1}_{h(x) \neq y} | X = x) &= \mathbb{P}(h(x) = 0 | X = x) \cdot \mathbb{P}(Y = 1 | X = x) + \mathbb{P}(h(x) = 1 | X = x) \cdot \mathbb{P}(Y = 0 | X = x) \\ &= \mathbb{P}(h(x) = 0 | X = x) \cdot (1 - p) + \mathbb{P}(h(x) = 1 | X = x) \cdot p \\ &\geq \underbrace{(\mathbb{P}(h(x) = 0 | X = x) + \mathbb{P}(h(x) = 1 | X = x))}_{1} \min(p, 1 - p) \end{aligned}$$

■

על כן $\forall h \in \mathcal{H} \quad L_D(f_D) \leq L_D(h)$ ומכאן האופטימליות כנדרש.

12.2 Naive Bayes

היות ואיננו יודעים את \mathcal{D} איננו יכולים להשתמש ב-*optimal Bayes*. קיימות מגוון שיטות המאפשרות לשערך צפיפות של דגימות ובכך לקבל מושג יותר טוב לגבי \mathcal{D} . כאשר מימד ה-*feature space* גבוה שיטות אלה יקרות ועל כן נניח "הנחה נאיבית" כי ה-*features* השונים בלתי תלויים. את כלל ה-*optimal Bayes* נוכל לכתוב בתור $h_D(x) = \underset{y=\pm 1}{\operatorname{argmax}} p_D(y|x)$ ובעזרת נוסחת בייס:

$$h_D(x) = \underset{y=\pm 1}{\operatorname{argmax}} \frac{p_D(x|y) p_D(y)}{p_D(x)} = \underset{y=\pm 1}{\operatorname{argmax}} p_D(x|y) p_D(y)$$

וכעת תחת ההנחה כי ה-*features* בלתי תלויים נקבל כי

$$p(x|y) = \prod_{i=1}^m p_D(x_i|y)$$

לכן מספיק שנדע את $p_D(x|y), p_D(y)$.

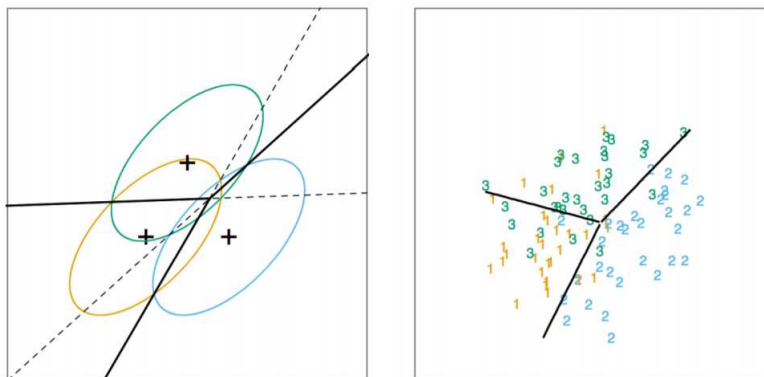
12.3 Linear Discriminant Analysis

כמו בדוגמה של *Naive Bayes* נרצה לתייג ל- $\{0, 1\}$ על בסיס $x = (x_1, \dots, x_d) \in \mathcal{X}$ כאשר הפעם ההנחה עבור \mathcal{D} היא $p_D(x|y) = \mathcal{D}[x|Y=y] = \mathcal{N}(\mu_y, \Sigma)$ אשר \mathcal{N} משתנה מקרי מעל \mathbb{R}^d המתפלג נורמלי עם תוחלת μ_y ושונות Σ .

הערה 12.2 פונקציית הצפיפות של משתנה שכזה היא:

$$\mathbb{P}(X = x | Y = y) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

כלומר אנחנו מניחים כי ההתפלגות של x בהינתן התפלגות התגיות היא גאוסיאן. על כן כאשר הגאוסיאנים רחוקים אחד מהשני הסיווג פשוט. כאשר הם קרובים ויש חפיפה הסיכוי לטעות גדל. קווי הגובה של פונקציה זו הם אליפסואידים והיות ומטריצת השונות זהה בין תיוגים שונים "אזורי התיוג" לכל תיוג (השטח במרחב שבתוכו נקבל את התיוג) הם אליפסואידים שווים בצורתם כאשר ההבדל הוא מיקום המרכז.



טענה 12.3 בהינתן D המקיימת את ההנחה לעיל ונסמן $D(Y = k) = \pi_k$ אזי ה-*bayes optimal* מתקבל על ידי $\operatorname{argmax}_{y=\pm 1} \delta_y(x)$ כאשר:

$$\delta_y(x) = x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu + \log(\pi_y)$$

הוכחה: תהי $x \in \mathbb{R}^d$ אזי

$$\begin{aligned} h_D(x) &= \operatorname{argmax}_y D(x|y) D(y) \\ &= \operatorname{argmax}_y \left(\det(2\pi\Sigma)^{-0.5} \pi_y \exp\left(\frac{1}{2}(x - \mu_y)^T \Sigma^{-1} (x - \mu_y)\right) \right) \\ &= \operatorname{argmax}_y \left(\log(\pi_y) + \frac{1}{2}(x - \mu_y)^T \Sigma^{-1} (x - \mu_y) \right) \\ &= \operatorname{argmax}_y \left(\log(\pi_y) + x^T \Sigma^{-1} \mu_y - \frac{1}{2} \mu_y^T \Sigma^{-1} \mu_y \right) \\ &= \operatorname{argmax}_y (\log(\pi_y) + \delta_y(x)) \end{aligned}$$

■

כעת כדי לממש זאת בהינתן $S = \{(x_i, y_i)\}_{i=1}^m$ נצטרך להעריך את μ_y, π_y, Σ . נסמן $m_y = \#\{Y = y\}$ (מספר הדוגמאות

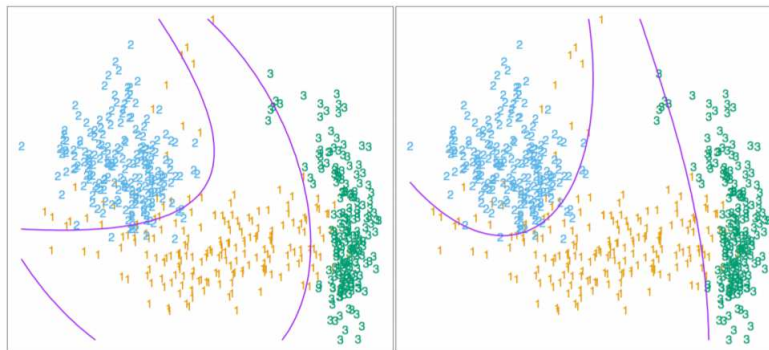
המתוייגות y עבור $y \in \mathcal{Y}$), $\hat{\pi}_y = \frac{m_y}{m}$, $\hat{\mu}_y = \frac{\sum_{i: y_i=y} x_i}{m_y}$, ולבסוף

$$\hat{\Sigma} = \frac{1}{m-2} \sum_{j=\pm 1} \sum_{i: y_i=j} \left(\sum (x_i - \hat{\mu}_j) (x_i - \hat{\mu}_j)^T \right)$$

נשים לב כי נוכל להציג כעת את δ בתור $\delta_y(x) = a_y + b_y^T x$ ומכאן כי גבול הסיווג הינו ליניארי ב- x . כלומר ה-*Bayes optimal predictor* במקרה זה הינו *halfspace*.

12.4 Quadratic Discriminant Analysis

ב-*LDA* הנחנו כי לכל תיוג מטריצת השונות בהתפלגות הנורמלית זהה. כעת נניח כי $\forall y \in \mathcal{Y}$ מטריצת השונות היא Σ_y (מקרה פרטי הוא ה-*LDA* שבו זהה). במקרה זה קיבלנו את ה-*QDA*.



Boosting 13

Boosting הינה פרדיגמה אלגוריתמית שבאה לעזור עם שתי בעיות. הראשונה היא ה-*bias – complexity tradeoff* והשנייה היא ה-*computational complexity* של הלמידה. ראינו תחת ב-*bias – complexity tradeoff* כי שגיאת המסווג ניתנת לפירוק ל-*approximation error* ול-*estimation error* וכי כיצד כאשר מחלקת היפותזות עשירה יותר אז חל שיפור ב-*approximation error* אבל קיימת עליה ב-*estimation error*. ב-*boosting* נרצה לקבל מעט יותר שליטה על *tradeoff* זה. במקרה של *computational complexity* נראה כי מימוש עיקרון ה-*ERM* במקרים רבים הינו בעייתי חישובית. שימוש ב-*boosting* מאפשר לשמר ביצועים של *weak learners* אשר נגדיר בפרק זה.

Weak Learnability 13.1

הגדרה 13.1 מחלקת היפותזות \mathcal{H} הינה (ϵ, δ) -weak-learnable אם קיים אלגוריתם למידה \mathcal{A} וסט אימון מגודל m כך שלכל הסתברות \mathcal{D} מעל \mathcal{X} ולכל $f \in \mathcal{H}$ מתקיים:

$$\mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}, f}(\mathcal{A}(S)) \leq \epsilon) \geq 1 - \delta$$

הערה 13.2 נשים לב כי בשונה מהגדרת ה-*PAC* (שנקרא לצורת למידה זו גם למידה חזקה, *strong learning*) ב-*weak learning* על האלגוריתם לספק את התנאים ל- ϵ, δ ספציפיים ולא $\forall \epsilon, \delta \in (0, 1)$.

הערה 13.3 עבור $\delta = 0$ למעשה עלינו לספק אלגוריתם שבהסתברות 1 בעל דיוק ϵ ועבור $\delta = 1$ תמיד נעמוד בתנאי. כאשר $\epsilon = 1$ תמיד נצליח לספק אלגוריתם שבכל הסתברות עומד בתנאי הדיוק וכאשר $\epsilon = 0$ למעשה נצטרך אלגוריתם עם איבוד אפס וראינו כי לא קיים כזה. לבסוף עבור $\delta = 0, \epsilon = \frac{1}{2}$, מדובר ב-*random* ולכן נניח לומד חדש בהכרח מקיים $\epsilon < \frac{1}{2}$.

דוגמה-מנחה: תהי $\mathcal{B} = \{x \mapsto \text{sign}(x - \theta) \cdot b : \theta \in \mathbb{R}, b \in \{\pm 1\}\}$ מחלקת ההיפותזות של *Decision Stamps* (באופן כללי $\{x \in \mathbb{R}^d \mapsto \text{sign}(x_j - \theta) \cdot b : \theta \in \mathbb{R}, b \in \{\pm 1\}, j \in [d]\}$ *decision stamps* הינם עץ החלטה בעל רמה אחת

$$\begin{array}{c} + \qquad \qquad - \\ \hline \theta \end{array}$$

כעת נגדיר $\mathcal{X} = \mathbb{R}$ ומחלקת היפותזות \mathcal{H} להיות מחלקת ה-*3-piece classifiers*. קל לראות כי לא נוכל ללמוד *PAC* את המחלקה הזאת בעזרת \mathcal{B} אולם נראה כיצד נוכל ללמוד אותה *weak-learning* בעזרת \mathcal{B} . מחלקה זו היא למעשה סיווג לפי שלושה חלקים:

$$\begin{array}{c} + \qquad \qquad - \qquad \qquad + \\ \hline \theta_1 \qquad \qquad \theta_2 \end{array}$$

טענה: קיים $m \in \mathbb{N}$ כך שעבור אלגוריתם המממש ERM_B האלגוריתם הינו $weak - learner$ $(\frac{5}{12}, \frac{1}{2})$ עבור \mathcal{H} .

הוכחה: נתחיל מטענת עזר:

למה: לכל \mathcal{D} מעל \mathcal{X} ולכל $f \in \mathcal{H}$ קיימת $h \in \mathcal{B}$ המקיימת $L_{\mathcal{D},f}(f) \leq \frac{1}{3}$.

הוכחה: עבור $x \in \mathcal{X}$ נסמן $\mathbb{P}(x \leq \theta_1) = \mathcal{D}_1, \mathbb{P}(\theta_1 < x \leq \theta_2) = \mathcal{D}_2, \mathbb{P}(x > \theta_2) = \mathcal{D}_3$. מתקיים $\mathcal{D}_1 + \mathcal{D}_2 + \mathcal{D}_3 = 1$ וכן קיים $i \in [3]$ כך ש- $\mathcal{D}_i \leq \frac{1}{3}$. נניח בלי הגבלת הכלליות כי $i = 1$ אזי $\mathcal{D}_1 \geq \frac{2}{3}$ ונבחר $\theta_1 = \theta_2 = x$.

$$h(x) = \begin{cases} + & x > \theta_2 \\ - & x \leq \theta_2 \end{cases}$$

ולמעשה האלגוריתם טועה רק על \mathcal{D}_1 ולכן: $L_{\mathcal{D},f}(f) = \mathcal{D}_1 \leq \frac{1}{3}$ כנדרש.

כעת היות ואנו יודעים את ה- $VCdim(\mathcal{B})$ ובפרט סופי נוכל, בעזרת המשפט המרכזי של הלמידה הסטטיסטית עבור $\delta = \frac{1}{2}$ נקבל שגיאיה של $\frac{1}{3} + \frac{1}{12} = \frac{5}{12}$ אשר נבחר בתור ϵ (נשים לב כי $\frac{5}{12} < \frac{1}{2}$) ועל כן קיים אלגוריתם $weak - learner$ $(\frac{5}{12}, \frac{1}{2})$ עבור \mathcal{H} כנדרש.

13.2 Boosting Confidence

יהי \mathcal{A} אלגוריתם $weak - learner$ (ϵ_0, δ_0) למחלקת היפותזות \mathcal{H} הדורש m_0 דגימות. נראה כיצד נוכל ללמוד את \mathcal{H} בדיוק של $\epsilon_0 + \epsilon$ וביטחון δ לכל $\epsilon, \delta \in (0, 1)$:

1. נריץ את \mathcal{A} על $k = \left\lceil \frac{\log(2/\delta)}{\log(1/\delta_0)} \right\rceil$ סטים $i.i.d$ כל אחד בגודל m_0 לקבלת h_1, \dots, h_k .

2. ניקח סט ולידציה V כך ש: $|V| \geq \frac{2\log(4k/\delta)}{\epsilon^2}$ ונחזיר $\hat{h} \in \underset{h_i}{\operatorname{argmin}} L_V(h_i)$.

טענה 13.4 $\mathbb{P}(L_{\mathcal{D}}(\hat{h}) \leq \epsilon_0 + \epsilon) \geq 1 - \delta$

הוכחה: תחילה נחסום את ההסתברות לשגיאה גדולה מ- ϵ_0 :

$$\begin{aligned} \mathbb{P}\left[\min_i L_{\mathcal{D}}(h_i) > \epsilon_0\right] &= \mathbb{P}[\forall i L_{\mathcal{D}}(h_i) > \epsilon_0] \\ &\stackrel{(i.i.d)}{=} \prod_{i=1}^k \mathbb{P}[L_{\mathcal{D}}(h_i) > \epsilon_0] \\ &\stackrel{weak-learner}{\leq} (\delta_0)^k \\ &\stackrel{choosing k}{\leq} \delta/2 \end{aligned}$$

כלומר $\mathbb{P}\left[\min_i L_{\mathcal{D}}(h_i) \leq \epsilon_0\right] \geq 1 - \frac{\delta}{2}$. היות איננו יודעים את \mathcal{D} אנו משתמשים ב- L_V כמדד שיערוך ל- $L_{\mathcal{D}}$. תהי פונקציית האיבוד ℓ בעלת טווח של $[0, 1]$ אזי מאי שיוויון הופדינג נקבל כי בהסתברות לפחות $1 - \delta$:

$$|L_V(h) - L_{\mathcal{D}}(h)| \leq \sqrt{\frac{\log(2/\delta)}{2m_0}}$$

מתהליך הולידציה אנו מקבלים

$$\mathbb{P}\left[L_{\mathcal{D}}(\hat{h}) > \min_i L_{\mathcal{D}}(h_i) + \epsilon\right] \leq \delta/2$$

כעת באמצעות חסם האיחוד נקבל את מה שהיה להוכיח.

תכונות הפרוצדורה:

1. נדרשנו למספר רב של דגימות: $m_0 \cdot \left\lceil \frac{\log(2/\delta)}{\log(1/\delta_0)} \right\rceil + \frac{2\log(4k/\delta)}{\epsilon^2}$
2. זמן הריצה הינו: $\left\lceil \frac{\log(2/\delta)}{\log(1/\delta_0)} \right\rceil \cdot \left(O(WL) + \frac{2\log(4k/\delta)}{\epsilon^2} \right)$ כאשר WL הינו האלגוריתם הלומד החלש.
3. הגדלנו את הביטחון אך שילמנו בדיוק: במקום ϵ_0 כעת $\epsilon_0 + \epsilon$.

13.2.1 Boosting Confidence בתוחלת

נשים לב כי בהינתן לומד \mathcal{A} המקיים:

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}} + \epsilon$$

נסמן $\theta = L_{\mathcal{D}}(\mathcal{A}(S)) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}$ אזי $\mathbb{E}_{S \sim \mathcal{D}^m} [\theta] \leq \epsilon$. θ הינו משתנה מקרי אי שלילי אזי מאי שיוויון מרקוב:

$$\mathbb{P}[\theta \geq 2\epsilon] \leq \frac{\mathbb{E}[\theta]}{2\epsilon} \leq \frac{1}{2}$$

וקיבלנו כי \mathcal{A} הינו $weak - learner$ $(2\epsilon, \frac{1}{2})$ עליו נוכל לעשות $boosting$ ולשפרו.

13.3 Boosting Accuracy באמצעות AdaBoost

עבור $boosting accuracy$ נראה את אלגוריתם $Adaptive boosting$ ($AdaBoost$) אשר מספק היפותזה עם $empirical risk$ נמוך. בהינתן קלט $S = \{(x_i, y_i)\}_{i=1}^m$ כך ש- $y_i = f(x_i)$ עבור $f \in \mathcal{H}$, כלשהי, האלגוריתם רץ מספר איטרציות כאשר בכל איטרציה מגדיר התפלגות $\mathcal{D}^{(t)}$ מעל S , מעביר את ההתפלגות WL ונותן ל- h_t שחזרה מה- WL משקל פרופורציונלי לשגיאה. באופן אינטואיטיבי השיטה מאלצת את הלומד להתמקד על הדוגמאות ה"קשות" יותר.

AdaBoost

input:
training set $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$
weak learner WL
number of rounds T
initialize $\mathbf{D}^{(1)} = (\frac{1}{m}, \dots, \frac{1}{m})$.
for $t = 1, \dots, T$:
 invoke weak learner $h_t = WL(\mathbf{D}^{(t)}, S)$
 compute $\epsilon_t = \sum_{i=1}^m D_i^{(t)} \mathbb{1}_{[y_i \neq h_t(\mathbf{x}_i)]}$
 let $w_t = \frac{1}{2} \log \left(\frac{1}{\epsilon_t} - 1 \right)$
 update $D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-w_t y_i h_t(\mathbf{x}_i))}{\sum_{j=1}^m D_j^{(t)} \exp(-w_t y_j h_t(\mathbf{x}_j))}$ for all $i = 1, \dots, m$
output the hypothesis $h_s(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^T w_t h_t(\mathbf{x}) \right)$.

כלומר, בהינתן כקלט סט אימון מתוייג, אלגוריתם למידה חלש ו- T מספר איטרציות האלגוריתם:

$$(0) \quad \text{מאתחל התפלגות אחידה } D^{(1)} \sim \text{Unif}([m]) \text{ מעל } S.$$

$$(1) \quad \text{לכל איטרציה } t:$$

(1.1) האלגוריתם קורא ל- WL על סט האימון S והתפלגות $\mathcal{D}^{(t)}$. ה- WL מחזיר היפותזה מתאימה h_t עבור ההתפלגות וסט האימון. נבחין כי מהגדרה של WL האלגוריתם, עבור ϵ_0, δ_0 מחזיר היפותזה מתאימה לכל התפלגות, ולכן בפרט עבור $\mathcal{D}^{(t)}$.

(1.2) נחשב את שגיאת ההכללה $\epsilon_t = L_{\mathcal{D}^{(t)}}(h_t) = \sum_{i=1}^m D_i^{(t)} \mathbb{1}_{y_i \neq h_t(x_i)}$ - כלומר מספר הדוגמאות שה- WL טעה עליהם עבור התפלגות $\mathcal{D}^{(t)}$.

(1.3) נגדיר $w_t = \frac{1}{2} \log \left(\frac{1}{\epsilon_t} - 1 \right)$ אשר ישמש כמשקולת עבור h_t בבניית ההתפלגות הבאה. נבחין כי בהינתן WL לומד חלש אנו מניחים כי $\epsilon_t < \frac{1}{2}$ (שכן אחרת לא שיפרנו מעל $random$) ועל כן w_t מספר חיובי. המשקולת הינה פרופורציונלית הפוכה לשגיאה - דוגמאות עליהן הצליח המסווג יקבלו משקל נמוך יותר ודוגמאות עליהן טעה המסווג יקבלו משקל גבוה יותר. כך נכריח את המסווג, באיטרציה הבאה, להצליח טוב יותר על הטעויות.

(1.4) נבנה התפלגות חדשה $\mathcal{D}^{(t+1)}$ באופן הבא: לכל דגימה $i \in [m]$ נגדיר את התפלגותה על ידי $\mathcal{D}_i^{(t+1)} = \frac{\mathcal{D}_i^{(t)} \exp(-w_t y_i h_t(x_i))}{\sum_j \mathcal{D}_j^{(t)} \exp(-w_t y_j h_t(x_j))}$. כאשר ההיפותזה מצליחה על דגימה מסוימת אזי $y_i h_t(x_i) > 0$ ועל כן היות ומוכפל ב- w_t - נקבל כי משקל הדגימה יקטן. כאשר ההיפותזה טועה על דגימה מסוימת אזי $y_i h_t(x_i) < 0$ ולכן כאשר מוכפל ב- w_t - נקבל כי משקל הדגימה יגדל.

(2) לבסוף נחזיר את סימן הסיווג (התיוג ± 1) על בסיס סכום משוקלל של כל ההיפותזות כפול משקלן.

טענה 13.5 שגיאת היפותזה h_t ביחס להתפלגות $\mathcal{D}^{(t+1)}$ היא $\frac{1}{2}$

הוכחה:

$$\begin{aligned} \sum_{i=1}^m \mathcal{D}_i^{(t+1)} \mathbb{1}_{y_i \neq h_t(x_i)} &= \frac{\sum_{i=1}^m \mathcal{D}_i^{(t)} e^{-w_t y_i h_t(x_i)} \mathbb{1}_{y_i \neq h_t(x_i)}}{\sum_{j=1}^m \mathcal{D}_j^{(t)} e^{-w_t y_j h_t(x_j)}} = \frac{e^{w_t \epsilon_t}}{e^{w_t \epsilon_t} + e^{-w_t (1-\epsilon_t)}} \\ &= \frac{\epsilon_t}{\epsilon_t + e^{-2w_t (1-\epsilon_t)}} = \frac{\epsilon_t}{\epsilon_t + \frac{\epsilon_t}{1-\epsilon_t} (1-\epsilon_t)} = \frac{1}{2} \end{aligned}$$

■

טענה 13.6 יהי WL אלגוריתם $weak - learner$ $(\frac{1}{2} - \gamma, \delta)$ אזי $\mathbb{P}(L_S(h_s) \leq \exp(-2\gamma^2 T)) \geq 1 - \delta T$

הערה 13.7 לכל $\epsilon > 0$ ו- $\gamma \in (0, \frac{1}{2})$ אם $T \geq \frac{\log(1/\epsilon)}{w\gamma^2}$ אזי $AdaBoost$ יחזיר היפותזה h_s עם $L_S(h_s) \leq \epsilon$

הערה 13.8 זמן הריצה של האלגוריתם הינו $T(O(m) + O(WL))$ - בכל איטרציה $t \in [T]$ האלגוריתם מריץ את ה- WL על S מחשב את ϵ_t על סמך S ($O(m)$), מחשב את w_t ($O(1)$) ומעדכן את ההתפלגות ($O(m)$).

דוגמה-מנחה כפי שראינו אלגוריתם ה- $AdaBoost$ למעשה מחזיר הרכבה (*ensemble/composition*) של מסווגים חלשים. תהי מחלקת ההיפותזות $\mathcal{B} = \{x \mapsto \text{sign}(x - \theta) \cdot b : \theta \in \mathbb{R}, b \in \{\pm 1\}\}$ של $Decision Stamps$ (שראינו בתחילת הפרק) אשר WL יחזיר היפותזה ממחלקה זו. אזי אלגוריתם $AdaBoost$ מחזיר היפותזה מהמחלקה הבאה:

$$L(\mathcal{B}, T) = \left\{ x \mapsto \text{sign} \left(\sum_{t=1}^T w_t h_t(x) \right) \mid w \in \mathbb{R}^T, \forall t \ h_t \in \mathcal{B} \right\}$$

כלומר כל $h \in L(\mathcal{B}, T)$ מכיל T היפותזות בסיסיות ב- \mathcal{B} הממושקלות לפי וקטור $w \in \mathbb{R}^T$. הפרדיקציה של h על דגימה x למעשה קודם יוצרת וקטור

$$\psi(x) = (h_1(x), \dots, h_T(x)) \in \mathbb{R}^T$$

אשר ה- $halfspace$ שמגדיר w מפריד אותו. על כן נוכל להגדיר את מחלקת ההיפותזות בתור:

$$L(\mathcal{B}, T) = \{x \mapsto \text{sign}(\langle w, \psi(x) \rangle) \mid w \in \mathbb{R}^T, \|w\|_0 \leq T\}, \quad \|w\|_0 = |\{i : w_i \neq 0\}|$$

13.3.1 *Bias – complexity tradeoff*

טענה 13.9 $VCdim(L(\mathcal{B}, T)) \leq \tilde{O}(T \cdot VCdim(\mathcal{B}))$

מסקנה 13.10 *AdaBoost* יחזיר היפותזה h_s אשר בהסתברות של לפחות $1 - \delta$: \tilde{O} - עד כדי פקטורים לוגריתמים)

$$L_{\mathcal{D}}(h_s) \leq L_S(h_s) + \tilde{O}\left(\sqrt{\frac{T \cdot VCdim(\mathcal{B}) + \log(1/\delta)}{m}}\right)$$

נובע מכך:

- ה-*estimation error* של *AdaBoost* גדל לינארית ב- T .
 - ה-*empirical risk* של *AdaBoost* קטן עם T .
 - באופן כללי ניתן להשתמש ב- T להקטנת ה-*approximation error* של $L(\mathcal{B}, T)$ ועל כן T יכול לשמש כמנוף ב-*tradeoff* בין *bias* ו-*complexity*.
- להסבר יותר מעמיק ניתן לקרוא בספר הקורס בעמוד 137.

14 Model Selection and Validation

כאשר ראינו בפרק הקודם את אלגוריתם *AdaBoost* (וכפי שנראה בדוגמאות נוספות תחת *regularization, decision trees, Soft – SVM*) ראינו כיצד פרמטר T שולט ב-*bias – complexity tradeoff*. כיצד נבחר T זה מבין $\{\mathcal{H}_T : T \in \mathbb{N}\}$?

14.1 Validation and k – Cross Validation

בתהליך ה-*Validation* נרצה להצליח לשערך כמה טוב המסווג שלנו. שיטה אחת, כפי שעשינו ב-*AdaBoost* (*Hold out set*), היא להחזיק סט דגימות נוסף לולידציה על המסווגים שקיבלנו ובעזרתו לבחור את האופטימלי.

אחת השיטות הנפוצות כדי להעריך *prediction error* הינה *Cross – Validation*. באופן טבעי, אם היה לנו מספיק *data* היינו שומרים חלק בצד בתור *validation set* כדי לבדוק את טיב המודל שנבחר. היות ו-*data* איננו בלתי מוגבל נרצה להשתמש במה שיש גם ל-*train* וגם ל-*validation*. עבור $K \in \mathbb{N}$ כאשר נדבר על *K – fold cross – validation* נחלק את הדגימות ל- K חלקים שווים באופן מקרי. נאמן מודל על סמך $K - 1$ מקטעים ונבדוק על המקטע ה- K . כל פעם נבחר את המקטע בדיקה (ולידציה) להיות מקטע אחר ולבסוף נמצע את כל התוצאות.

1	2	3	4	5
Train	Train	Validation	Train	Train

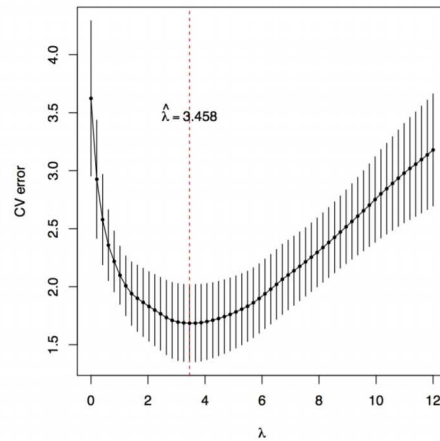
באופן יותר פורמלי תהי $k : [N] \rightarrow [K]$ הממפה דגימה ל-*partition* באופן מקרי. נסמן $\hat{f}^{-k}(x)$ ההיפוזתה שנבחרה על ה-*data* ללא ה-*partition* ה- k . אזי שיערוך השגיאה באמצעות ה-*cross validation* הינו:

$$CV(\hat{f}) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{-k(i)}(x_i))$$

אם אנו עוסקים במודל בעל פרמטר כיוול α כמו למשל פרמטר הרגולריזציה בעצי החלטה או *Soft SVM* או כדוגמת פרמטר k ב-*KNN* (ראו בהמשך) אז אנו למעשה עוסקים במשפחה של מחלקות היפותזות \mathcal{H}_α . נסמן $\hat{f}^{-k}(x, \alpha)$ המודל ה- α שנבחר עבור הפרמטר רגולריזציה α שנבחר על ה-*data* ללא ה-*partition* ה- k . אזי שיערוך השגיאה באמצעות ה-*cross validation* הינו:

$$CV(\hat{f}, \alpha) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{-k(i)}(x_i, \alpha))$$

במקרה שכזה נוכל לשערך את טעות ה- CV כתלות בפרמטר α (בגרף) עבור ערכים שונים ולקבל עקומה שתעזור לנו לקבוע מה ה- λ המתאים לבעיה.



Pseudocode

```

 $k$ -Fold Cross Validation for Model Selection

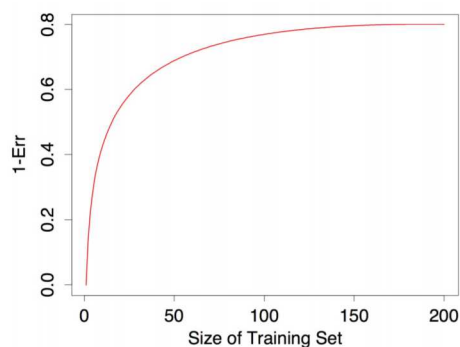
input:
  training set  $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ 
  set of parameter values  $\Theta$ 
  learning algorithm  $A$ 
  integer  $k$ 
partition  $S$  into  $S_1, S_2, \dots, S_k$ 
foreach  $\theta \in \Theta$ 
  for  $i = 1 \dots k$ 
     $h_{i,\theta} = A(S \setminus S_i; \theta)$ 
     $\text{error}(\theta) = \frac{1}{k} \sum_{i=1}^k L_{S_i}(h_{i,\theta})$ 
output
   $\theta^* = \text{argmin}_{\theta} [\text{error}(\theta)]$ 
   $h_{\theta^*} = A(S; \theta^*)$ 

```

14.1.1 בחירת ערך ה- k עבור k -fold

כאשר $k = 1$ כל ה- $data$ הזמין משמש לאימון ולמעשה אין $cross\ validation$. כאשר $k = 2$ (המכונה גם $split - sample$) רק מחצית המידע משמש ל- $train$ ולמעשה השגיאה שנקבל עם CV תהיה גדולה מבלי. באופן כללי עבור k קטנים אנחנו עלולים להתאמן על $datasets$ קטנים מדי ונקבל כי ה- $CV\ error$ הוא $biased\ upwards$ - משערכים שגיאה גדולה ממה שנוכל להשיג בפועל. עבור ערכי k גדולים ההבדל בין קבוצות המדגם באימונים השונים נמוך. על כן תהיה קורלציה גבוהה יותר בין תוצאות אימונים שונים ולמעשה נגדיל את ה- $variance$. ככלל אצבע, על מספיק $data$, $k = 5, 10$ נותנים ביצועים טובים (ניתן לקרוא את *Breiman and Spector, 2012* או *Kohavi, 1995*).

היות ובחירת k משפיעה גם על גודל סט האימון נרצה לראות עבור הדומיין הספציפי מהו גודל סט האימון הנחוץ. נוכל לצייר עקומה שתעיד על ההצלחה כתלות בגודל ה- $trainingset$. נראה כי עבור $N = 200$ או $N = 150$ ההצלחה (או 1 פחות הטעות) כמעט לא משנה ועל כן נוכל לבחור k שבו נאפשר פחות דגימות בסט האימון, נגדיל את השונות בין הסטים, נקטין את ה- $variance$ וזאת מבלי לאבד כמעט באיכות המסווג שיבחר.



14.1.2 איך לא בוחרים k

נניח כי אנו חוקרים בעיית למידה מסויימת ומנסים למצוא מודל מתאים. אנחנו מחפשים סט $predictors$ טובים במחלקה שלנו (ונניח כי המחלקה עשירה). לאחר שמצאנו סט שכזה המראה קורלציה טובה עם התיוגים נמצא מסווג על בסיסן. כעת נשתמש ב- $cross validation$ כדי להעריך את ערך פרמטר הכיול ולהעריך את ה- $prediction error$ הסופי.

מקרה זה מתאר מקרה שבו הכנסנו אופטימיות למודל. סט ההיפותזות שהתבססנו עליהן נבחר על סמך ביצועיהן על כל ה- $data$ ועל כן התוצאות שקיבלנו בסופו של תהליך (ה- $CV error$) הוא $biased$ לטובת היפותזות אלה ולמעשה אנחנו מנבאים טעות נמוכה מדי. הדרך הנכונה לפעול היא:

1. נפצל את הדגימות ל- K קבוצות ($folds$) באופן רנדומלי.

2. לכל $k = 1, \dots, K$:

(א) נמצא תת קבוצה של $predictors$ בעלות קורלציה טובה בין תוצאותיהן לתיוגים. בשלב זה נשתמש בכל ה- $samples$ למעט אלה ב- k -fold.

(ב) נמצא מסווג טוב על בסיס תת קבוצה זו בהסתכל על כל ה- $samples$ למעט אלה ב- k -fold.

(ג) נבדוק את טיב המסווג על הדגימות ב- k -fold.

Train – Validation – Test split **14.2**

בדרך כלל אלגוריתמים הכוללים *Validation* הם מהצורה:

1. *Training Stage* - מציאת ההיפותזה האופטימלית מכל מודל:

For each $\theta \in \Theta$: *Draw new* S

$$h_\theta = \mathcal{A}(S, \theta)$$

2. *Validation Stage* - מציאת היפותזה אופטימלית מעל כל סט המודלים:

Sample new set V

Choose $h^* = h_{\theta^*}$: $\theta^* = \underset{\theta \in \Theta}{\operatorname{argmin}} L_V(h_\theta)$

3. *Test Stage* - שיערוך טיב ההיפותזה שנבחרה:

Sample new set T

Estimate error by $L_T(h^*)$

Regularization 15

כאשר אנו באים לענות על בעיית למידה וקיים לנו דומיין \mathcal{X} על $feature\ space$ כלשהו אנו מעוניינים למצוא היפותזה הפותרת טוב את הבעיה. ראינו ב- $bias - complexity\ tradeoff$ שנרצה למצוא היפותזה אשר מידת המורכבות שלה לא נמוכה מדי ולא גבוהה מדי. בפרק זה נראה כלים העוזרים לנו למצוא היפותזה שכזו. בנוסף בבעיות למידה רבות אנו נתקלים במצבים שבהם מספר ה- $features$ גבוה מאוד ($d \sim m$ או אפילו $d > m$). כאשר מספר ה- $features$ רב, היכולת ליצור מודלים מתאימים טוב יתר ל- $data$ (ומכאן גם מורכבים מאוד) היא רבה ועל כן נרצה להיות ערים וזהירים מכך. נרצה גם להשתמש רק ב- $features$ שעוזרים לנבא הכי טוב את התיוגים מתוך כלל ה- $features$, ולהתעלם מ- $features$ שמידת התרומה שלהם לפרדיקציה היא נמוכה. לזאת קוראים $feature\ selection$ ונתעסק גם בזאת מעט בפרק זה.

15.1 עיקרון RLM – Regularized Loss Minimization

ראינו כי בעיקרון ה- ERM נחזיר היפותזה אשר מביאה למינימום את ה- $empirical\ risk$. במקרים רבים שימוש ב- ERM פשוטו כמשמעו יביא לקבלת היפותזה מורכבת אשר עונה טוב מאוד על סט האימון (מצליחה לתפוס את הנקודות ולמעשה "הולכת אחרי הרעש") אך זו תצליח פחות טוב על סט הבדיקה. במקרה זה נסבול מ- $variance$ גבוה ו- $over\ fitting$. נרצה להתאים את כללי הלמידה באופן כזה שיגבילו את מורכבות ההיפותזה המוחזרת ($regularized$).

הגדרה 15.1 פונקציית רגולריזציה ($regularizer$) הינה מיפוי $\mathcal{R} : \mathcal{H} \rightarrow \mathbb{R}$ כאשר עיקרון ה- RLM מחזיר היפותזה המביאה למינימום את
$$\sum_{i=1}^m (L(y_i, h(x_i)) + \lambda \mathcal{R}(h))$$
 עבור $\lambda \geq 0$.

λ משמש בתור ה- $complexity\ parameter$. באופן אינטואיטיבי מורכבות ההיפותזה נמדדת על ידי הערך שמוחזר מ- \mathcal{R} והאלגוריתם אם כך יחפש איזון ($tradeoff$) באמצעות λ בין $empirical\ risk$ נמוך לבין מודלים פשוטים. עבור $\lambda = 0$ נקבל את עיקרון ה- ERM ועבור λ גדול, היות ונחפש מינימום, נקבל רק מודלים פשוטים. כעת נראה מספר מימושים אפשריים של עיקרון זה ונעמוד על ההבדלים ביניהם. במבט כללי ההבדל ביניהם הוא השימוש בנורמות שונות ($\|\cdot\|_0, \|\cdot\|_1, \|\cdot\|_2$) כאשר בפועל "נורמת 0" איננה באמת נורמה) בפונקציית הרגולריזציה. נתבסס על רגרסיה לינארית כדוגמה מנחה על מנת להבין את המימושים השונים.

15.1.1 Best Subset Regression

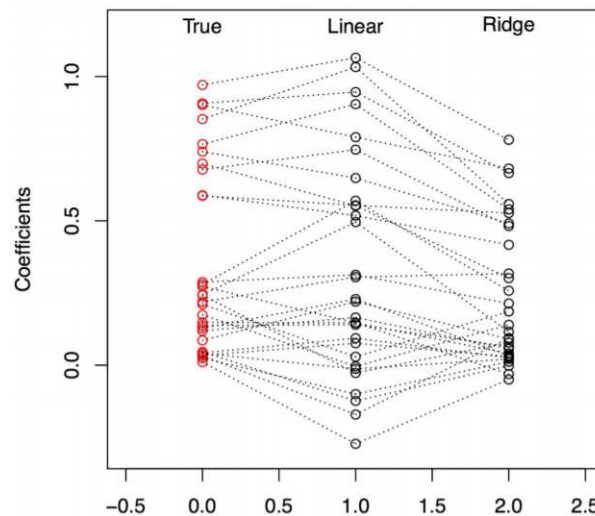
נגדיר $\mathcal{R}(w) = \#\{i : w_i \neq 0\}$ ונראה לפונקציה זו ה- $sparsity$ של w . נסמן גם ב- $\|w\|_0$. כלומר באופן אינטואיטיבי $regularizer$ זה יבקש לצמצם את מספר ה- $features$ אשר המודל משתמש בהם: נחפש מינימום גם על גורם זה - כלומר היפותזה המכילה פחות מקדמים שהם לא אפס - כלומר היפותזה שבה יש יותר $features$ שלא בשימוש. במקרה של רגרסיה לינארית שבה $\mathcal{H} = \{x \mapsto \langle x, w \rangle\}$ נקבל שנרצה מינימום של:

$$\|y - X^T w\|_2^2 + \lambda \|w\|_0$$

למרות שנראה כי הגדרה זו היא בדיוק מה שאנו מחפשים כדרך לומר כמה פשוט/מורכב נרצה את המודל שלנו, נשים לב שלמעשה ניתן לעשות רדוקציה מבעיה זו לבעיית מציאת תת קבוצה אופטימלית אשר הינה בעיה $NP - hard$. לכן לא קיים אלגוריתם יעיל לפתרון בעיה זו. בכל זאת, למספר d features נמוך (בערך עד כדי 40) קיימות ספריות שנותנות מענה.

Ridge Regression 15.1.2

ה-*Ridge Regression* מכווץ את המקדמים של הוקטור ברגרסיה על ידי אכיפת עלות לגודלם. בגרף הבא נתבונן בערכי המקדמים ל-*features* השונים על *data* שיצרנו ואנו יודעים מהם הערכים האמיתיים של המקדמים. נראה כי כאשר ביצענו את הרגרסיה הלינארית כל *feature* קיבל כעת מקדם כלשהו כאשר גם הפיזור (*variance*) גדולה יותר וגם כאשר נתבונן ב-*features* שהמקדמים שלהם הם 0 או כמעט 0 נראה כי חלקם קיבלו ערכים שונים מכך אשר למעשה פוגע בפרדיקציה. כאשר נשווה זאת עם תוצאת ה-*ridge regression* נראה כי קיבלנו הן פיזור צפוף יותר (*variance* נמוך יותר) וגם *features* עם מקדמים סביב ה-0 נשארו עם מקדמים שכאלה - כלומר כמו במקרה האמת הם אינם משמשים לצורך הפרדיקציה (ונוכל לשקול להורידם).



כאשר ניקח את *ERM* ועבור \mathcal{R} נבחר את $\|\cdot\|_2^2$ (ולעיתים מסמנים את הנורמות גם בתור ℓ_2 norm $\|\cdot\|$) נקבל את $\mathcal{A}(S) = \underset{w}{\operatorname{argmin}} \left(L_S(w) + \lambda \|w\|_2^2 \right)$: *Tikhonov regularization*. כשנשליך זאת על רגרסיה לינארית עם *squared loss* נקבל את ה-*Ridge regression*:

$$\|y - X^T w\|_2^2 + \lambda \|w\|_2^2$$

כשנססה לפתור בעיית מינימיזציה זו נקבל את המשוואות הנורמליות הבאות:

$$Xy = (XX^T + \lambda I) w \iff (XX^T + \lambda I) Xy = w$$

אשר $(XX^T + \lambda I)$ מטריצה הפיכה (כזכור XX^T איננה בהכרח הפיכה) - סימטרית חיובית עם ע"ע אי שליליים. על כן בדומה לרגרסיה לינארית קיים פתרון סגור. כפי שהגדרנו Σ^\dagger נגדיר:

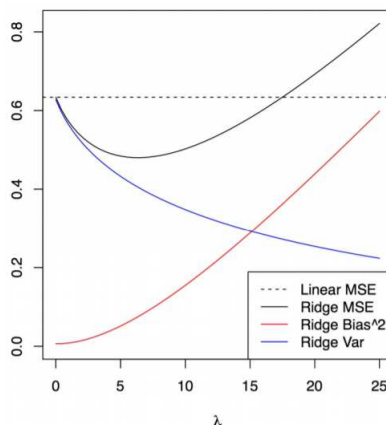
$$\Sigma_{i,i}^\lambda = \frac{\sigma_i}{\sigma_i^2 + \lambda}$$

ועבור פירוק ה- SVD של $X = U\Sigma V^T$ הפתרון הינו:

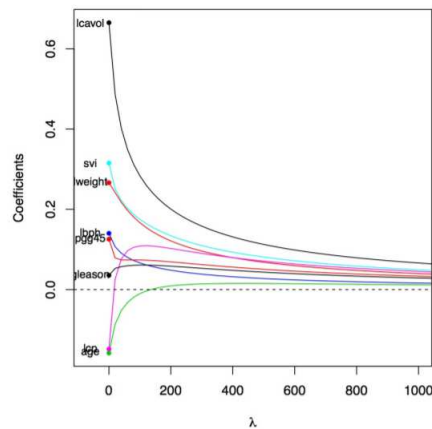
$$\hat{w}_\lambda = U\Sigma^\lambda V^T y$$

Bias variance tradeoff

נסתכל על ה-*bias variance tradeoff* כתלות ב- λ . עבור $\lambda = 0$ קיבלנו למעשה את ה- ERM ומודל מורכב יותר. עבור $\lambda = 25$ קיבלנו בהכרח מודל פשוט שכן מודלים מורכבים נקנסו ועל כן לא יבחרו. באופן כללי נראה כי ה-*variance* קטן בעוד שה-*bias* גדל. היות וככל ש- λ גדל אנו מכריחים מודלים פשוטים יותר היכולת שלהם ללכת אחרי הרעש יורדת והתוצאות בהתאם. נראה כי אכן ישנם ערכי λ עבורם ה-*ridge regression* מצליח להוריד *variance* בקצב מהיר יותר מעליית ה-*bias* - אלה הם ערכי ה- λ שנרצה לבחור.



Regulaization Path כאשר נקבע בעיה ספציפית ונבדוק את שינוי המקדמים של ה-*features* כתלות ב- λ נקבל את הגרף הבא. גרף זה מכונה ה-*regularization path*. נראה כי על אף ש-*Ridge* משפר את *MSE* וכאן כאשר $\lambda \rightarrow \infty$ המקדמים שואפים ל-0 הוא איננו נותן לנו אינטואיציה טובה איזה *features* נוכל להוריד ולא להשתמש בהם.

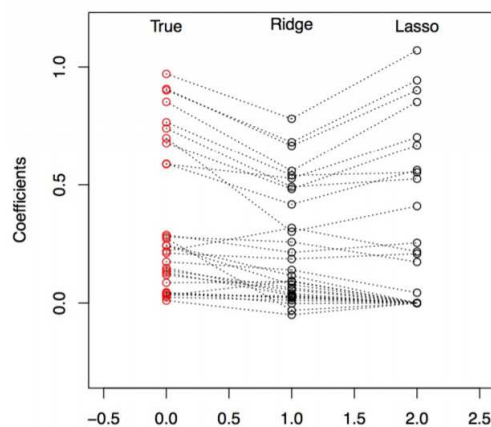


L_1 Lasso 15.1.3

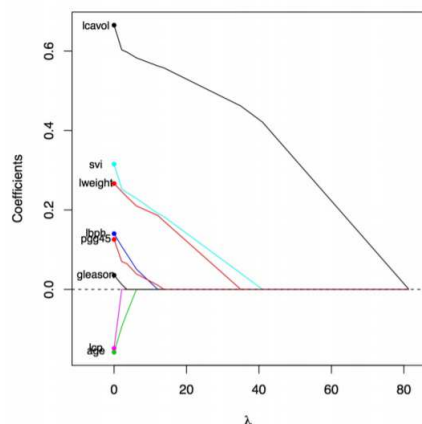
בדומה ל- $ridge$ ה- $lasso$ הינה פונקציית כיווץ מקדמים כאשר ה- $penalty$ הינו לפי $\|\cdot\|_1$. כלומר $\mathcal{R}(w) = \sum_i |w_i|$.

$$\operatorname{argmin}_w \|y - X^T w\|_2^2 + \lambda \|w\|_1$$

השימוש ב- $lasso$ מוביל לכך שההיפותזה המוחזרת אינם לינאריים ב- y_i ואין נוסחה סגורה למציאת ההיפותזה המינימלית. עם RLM (אשר הינו בעיית למידה קמורה - ראו בהמשך הסיכום) הינו $quadratic programming$ ולכן קיימים אלגוריתמים יעילים לפתרון בעיית אופטימיזציה. בדומה ל- $ridge$ כאשר $\lambda = 0$ נקבל את כלל ה- ERM וככל ש- λ גדל נחפש היפותזות אשר $\|w\|_1$ קטן אשר מוביל למעשה לפתרון שבו יותר ויותר מהמקדמים הם 0. מכאן כי ה- $constraint$ ש- $lasso$ גורם הוא למעשה סוג של $subset selection$. כאשר נבדוק את פיזור המקדמים עבור $lasso$, בדומה למה שעשינו עבור $ridge$ נראה כי בשונה מ- $ridge$ עבור $features$ שבקושי תרמו לפרדיקציה (מקדמים קרובים ל-0) ה- $lasso$ נתן מקדם 0. כלומר האלגוריתם החליט כך ש- $features$ אלא אינם רלוונטיים לפרדיקציה.



Regulaization Path כאשר נבדוק את ה-*regulaization path* עבור *lasso* נקבל את התמונה הבאה. נראה כי בשונה מהשינוי האקספוננציאלי שקיבלו המקדמים ב-*ridge*, במקרה הזה השינוי הינו כמעט ליניארי ב- λ . כלומר נוכל אפוא להחליט מה יהיה גודל *subset* של ה-*features* שנרצה (ובכך להגביל את ה-*complexity* של ההיפותזה הנבחרת) ואז בהתאם לבחור את ה- λ . לדוגמה עבור $\lambda = 20$ בדוגמה הזאת נקבל כי מחצית ה-*features* קיבלו מקדם 0 ולא ישחקו תפקיד בפרדיקציה.



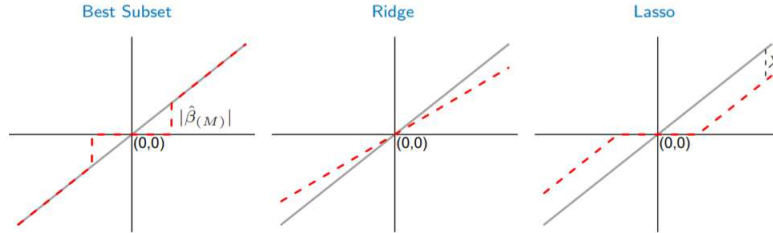
15.2 השוואה בין פונקציות הרגולריזציה

תחילה נתייחס ל-*intercept*. הן ב-*ridge* והן ב-*lasso* לא מבצעים רגולריזציה ל-*intercept*. כלומר בנורמה, בפעולת הסכימה עובדים רק על הקואורדינטות 1 עד d ולא מ-0. ספציפית למקרה של *lasso* במקרים רבים נתאים את המטריצה X כך ששורות המטריצה היו בעלות תוחלת אפס. כעת ננסה להבין ממה נובעים ההבדלים בין *ridge* ו-*lasso*.

15.2.1 Hard/Soft Thresholding

במקרה שבו X אורתונורמלית הפונקציות השונות משרות טרנספורמציה פשוטה על ה-*Least Squares*. נשים לב כי *ridge* לא שולח אף ערכים ישירות לאפס בעוד ש-*best subset* ו-*lasso* כן. בנוסף *Lasso* מכווץ מקדמים שאינם נשלחים ל-0 בעוד ש-*best subset* לא משתנים המקדמים של ה-*features* שלא נשלחו ל-0.

Estimator	Formula
Best subset (size M)	$\hat{\beta}_j \cdot I(\hat{\beta}_j \geq \hat{\beta}_{(M)})$
Ridge	$\hat{\beta}_j / (1 + \lambda)$
Lasso	$\text{sign}(\hat{\beta}_j)(\hat{\beta}_j - \lambda)_+$



ה-*best subset* נותן ערך 0 לכל מקדם שקטן מ- M המקדמים הראשונים. זה נקרא *hard thresholding*.
ridge regression מכווץ את המקדמים באופן פרופורציונלי ו-*lasso* מכווץ את המקדמים בפקטור λ ומאפס ערכים קטנים מ-0. שני אלה נקראים *soft thresholding*. כלומר עבור $\lambda \geq 0$ ו- $x \in \mathbb{R}$:

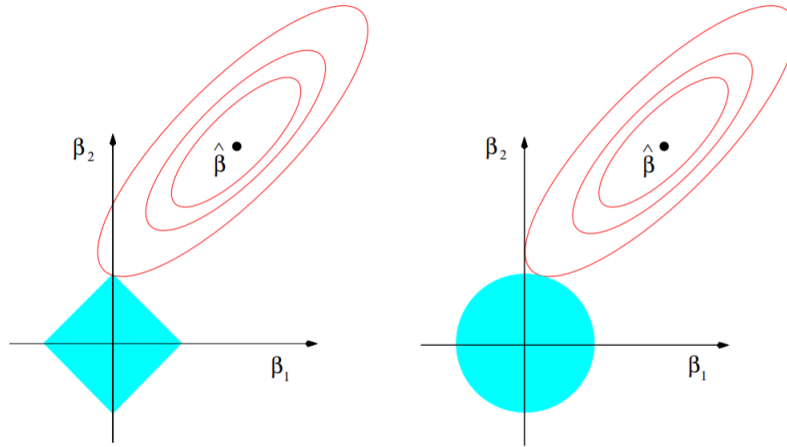
$$\eta_{\lambda}^{soft}(x) = \begin{cases} x - \lambda & x \geq \lambda \\ 0 & -\lambda > x > \lambda \\ x + \lambda & -\lambda \leq x \end{cases}, \quad \eta_{\lambda}^{hard}(x) = x \cdot \mathbb{1}_{|x| \geq \lambda}$$

ונקבל:

$$\begin{aligned} \hat{w}_{\lambda}^{ridge} &= \hat{w}^{LS} / (1 + \lambda) \\ \hat{w}_{\lambda}^{lasso} &= \eta_{\lambda}^{soft}(\hat{w}_{\lambda}^{LS}) \\ \hat{w}_{\lambda}^{subset} &= \eta_{\sqrt{\lambda}}^{hard}(\hat{w}_{\lambda}^{LS}) \end{aligned}$$

15.2.2 כדורי יחידה של הנורמות ו-*sparsity*

בדוגמה בה אנו עוסקים, פונקציית האיבוד *Least square*, עבור רגרסיה לינארית, הינה תבנית ריבועית ומכאן כי היא קמורה וקווי הגובה שלה אליפסואידים. הנקודות $\hat{\beta}$ מייצגות את נקודת המינימום של הפונקציות. כאשר נסתכל על הבעיה ב- \mathbb{R}^2 ונמקם זאת ליד כדורי היחידות של הנורמות השונות ($\|\cdot\|_1$ - משמאל, $\|\cdot\|_2$ - מימין) נקבל את נקודות החיתוך בין קווי הגובה השונים לבין כדורי היחידה. במקרה של $\|\cdot\|_1$ (שימוש ב-*lasso*) הנקודות השפיציות הן נקודות שבהן המקדם של *feature* מתאפס. על כן עיקרון ה-*RLM* יחזיר היפותזות אשר נמצאות בנקודות מפגש אלה. הקנס על שימוש בעוד *features* ידחוף לבחירת היפותזה על הקודקוד אשר שם יש פחות *features* עם מקדמים שונים מאפס. כאשר נדמיין זאת ב- \mathbb{R}^d נקבל "כדור" יחידה עם קודקודים רבים - כאשר בכל קודקוד כמות ושילוב שונה של *features* שהמקדמים שלהם הם 0.



15.3 *Lasso* עבור *Logistic Regression*

כשם ש- ℓ_1 regularization עובד טוב עבור רגרסיה לינארית הוא גם עובד טוב עבור רגרסיה לוגיסטית. כאשר ניקח את הרגרסיה לוגיסטית עם ה- $\log - \text{likelihood}$

$$\ell(\beta_0, \beta) = \sum_{i=1}^n y_i (\beta_0 + x_i^T \beta) - \sum_{i=1}^n \log(1 + \exp(\beta_0 + x_i^T \beta))$$

ננרמל ונוסיף *lasso* ל- $\log - \text{likelihood}$ negative נקבל את הפונקציית *loss* הבאה למצוא לה מינימום:

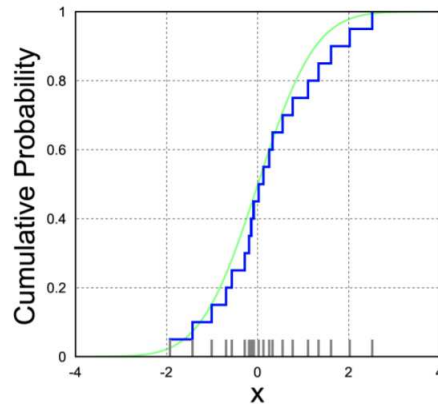
$$- \left[\frac{1}{m} \sum_{i=1}^m y_i (\beta_0 + x_i^T \beta) - \log(1 + \exp(\beta_0 + x_i^T \beta)) \right] + \lambda \|\beta\|_1$$

נשתמש ב-*cross validation* למצוא λ טוב (והיות ובעיה קמורה קיימים אלגוריתמים גרייס יעילים). למעשה באימפליקציה זו נקבל כי ℓ_1 יוצר לנו *sparse*-יות של וקטור המקדמים β - עבור λ גדול נקבל יותר מקדמים 0 - יותר *sparse* - פחות *features* שהינם חלק מהפרדיקציה. יוצא שמסווג זה הינו עם *misclassification rates* טובים ובוחר רק את ה-*feature subset* שרלוונטי לבעיה.

16 Bootstrapping and Bagging

16.1 Bootstrapping

Bootstrapping הינו כלי כללי לשיערוך דיוק סטטיסטי. נניח כי בידינו $z_1, \dots, z_n \stackrel{iid}{\sim} \mathcal{D}$ דגימות מעל התפלגות \mathcal{D} לא ידועה. נשים לב כי בהינתן n הדגימות הללו יש בידינו התפלגות אמפירית $\hat{\mathcal{D}}_n$ של דגימות אלו. כעת נניח כי בידינו איזושהי סטטיסטיקה (פונקציה על ה- data) $\hat{\theta}$ וברצוננו לשערך את $\hat{\theta}(X)$ כאשר $X \sim \mathcal{D}$ (כלומר מתוך ההתפלגות הכללית אשר איננה ידועה לנו). אם נניח כי $\hat{\mathcal{D}}_n \approx \mathcal{D}$ אזי עבור דגימה כרצוננו מהתפלגות של $\hat{\theta}(X^*)$ כאשר $X \sim \hat{\mathcal{D}}_n$ נקבל כי $\hat{\theta}(X^*) \approx \hat{\theta}(X)$. כעת נוכל לשלוף B דגימות באופן *monte – carlo* מ- $\hat{\theta}(X^*)$ ונקבל כי הדגימות דומות להתפלגות המקורית (והלא ידועה)



באופן ספציפי עבור בעיות למידה, עבור $B \in \mathbb{N}$ שנחליט, נבצע את התהליך הבא: נשלוף באופן מקרי, עם חזרות, m דגימות מתוך ה- dataset שנתון לנו. נסמן דגימות אלה בתור $(x_i^{(b)}, y_i^{(b)})_{i=1}^m$ ונכנה בשם *bootstrap samples*. עבור הסט ה- b נאמן מודל $h^{(b)}$. נוכל כעת למצע על סמך כל אוסף המסווגים הללו ולשערך את דיוקם טוב יותר.

16.2 Bagging

כאשר ראינו את *bootstrapping* המטרה הייתה לשערך את הדיוק שבפרדיקציה. כעת נשתמש ב-*bootstrapping* כדי לשפר את הפרדיקציה עצמה. יהי סט האימון $S = \{(x_i, y_i)\}_{i=1}^N$ ונניח לשם הדוגמה כי ביצענו רגרסיה לינארית ובידינו היפותזה $\hat{f}(x)$ ב-*Bagging* (או *Bootstrap aggregation*) נמצע את הפרדיקציה על קבוצת *bootstrap samples*. כלומר:

• עבור $B \in \mathbb{N}$ שנבחר, לכל $b \in [B]$:

- נשלוף m דגימות באופן מקרי עם חזרות $Z^{*b} = (x_i^{(b)}, y_i^{(b)})_{i=1}^m$.

- נאמן מודל חדש על הסט Z^{*b} שיספק את הפרדיקציה $\hat{f}^{*b}(x)$.

- ה-*estimation* עבור דגימה חדשה יהיה *ensemble* של המסווגים שיצרנו: (נוכל להחליט על שיטות מיצוע אחרות כמו *majority vote*)

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

באופן כללי כאשר $B \rightarrow \infty$ נקבל כי $\hat{f}_{bag}(x) \rightarrow \hat{f}(x)$. כלומר ה-*bagged estimation* שואפת לפרדיקציה אילו ידענו את \mathcal{D} (בספר הקורס Elements of Statistical Learning עמוד 282 יש הסבר תחת אילו תנאים הדבר איננו מתקיים). בשיטה זו נקטין את ה-*variance* מבלי להגדיל את ה-*bias*.

Bagging vs, Boosting

	Bagging	Boosting
sampling training data	uniform	non-uniform
Can we train in parallel?	yes	yes
weights over the classifiers	uniform	non-uniform
bias-variance trade-off	reduce variance	reduce bias

Feature Selection 17

בשלב זה ראינו כיצד מתנהג ה-*bias – complexity tradeoff*, כיצד לבחור מודל בהתאם לפרמטר כיוול (בהתאם ולמודל קיים) ושיטות שונות לשיפור יכולות הלימוד בהינתן *data* מוגבל. כאשר מימד הקלט (הן *feature space* והן כמות דגימות) גבוה קשה לנו להתמודד עם זאת. בהמשך נראה שני אלגוריתמים, *SVM* ו-*Kernels* שעוזרים להתמודד עם הבעיה. כרגע נדון כיצד נוכל לפשט בעיה זאת על ידי צמצום חלק מה-*features*. נשים לב כי ניתן להשיג מטרה זו באופן חלקי גם באמצעות שימוש ב-*Lasso* תחת *Regulaization* אולם זמן העבודה (אימון) אך גם עיבוד ה-*data* ואולי גם גורמים אחרים כמו עליות להשגתו) עדין קיימות. על כן נרצה למצוא איזה *features* נוכל להורידם באופן כללי. בספר ישנן מספר גישות לכך (עמוד 358) אך בקורס כיסינו רק את *Filter*.

בשיטה זו, שאולי מהפשוטה שבשיטות, אנו בוחנים *features* באופן בלתי תלוי מהשאר ומספקים "ניקוד" ל-*feature*. ערך זה יעיד על מידת הקשר בין התיוג של הדגימה לערך ה-*feature*. כמובן שגם לאופן חישוב הניקוד יתכנו פונקציות רבות. בחירה טבעית אחת תהיה ה-*Empirical Square Loss* של *predictor* לינארי:

$$score(j) = \min_{a,b \in \mathbb{R}} \frac{1}{m} \sum_{i=1}^m \left(a(x_i)_j + b - y_i \right)^2$$

כאשר $\{(x_i, y_i)\}_{i=1}^m \subset \mathbb{R}^d \times \mathbb{R}$ ו- $j \in [d]$ מייצג $feature$ ספציפי. פונקציה זו תכווין ל- $features$ בעלי ניקוד נמוך. נסמן $\bar{y} = \frac{1}{m} \sum y_i$ ו- $\bar{v} = \frac{1}{m} \sum v_i$ אזי עבור $v^j = ((x_1)_j, \dots, (x_m)_j)^T \in \mathbb{R}^m$, $1_m = (1, \dots, 1)^T \in \mathbb{R}^m$ נקבל:

$$\begin{aligned} score(j) &= \min_{a, b \in \mathbb{R}} \frac{1}{m} \|av^j + b \cdot 1_m - y\|_2^2 \\ &= \min_{a, b \in \mathbb{R}} \frac{1}{m} \|a(v^j - \bar{v}^j \cdot 1_m) + b \cdot 1_m - (y - \bar{y} \cdot 1_m)\|_2^2 \end{aligned}$$

וכדי למצוא a, b מתאימים נגזור לפי b ונשווה לאפס. נקבל כי $b = 0$. כעת נחשב עבור a ונקבל כי

$$a = \frac{\langle v^j - \bar{v}^j \cdot 1_m, y - \bar{y} \cdot 1_m \rangle}{\|v^j - \bar{v}^j \cdot 1_m\|_2^2}$$

כאשר נציב בפונקציה נקבל כי

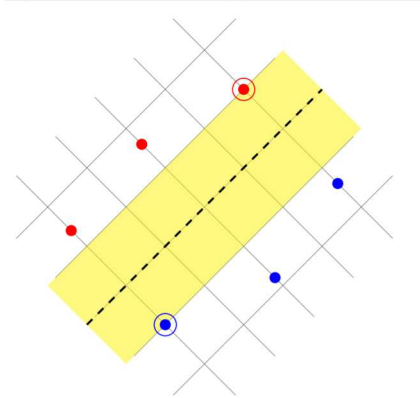
$$score(j) = \frac{1}{m} \left(\|y - 1_m \cdot \bar{y}\|_2^2 - \frac{\langle v^j - \bar{v}^j \cdot 1_m, y - \bar{y} \cdot 1_m \rangle^2}{\|v^j - \bar{v}^j \cdot 1_m\|_2^2} \right)$$

כעת נוכל על בסיס פונקציה זו להגדיר פונקציות $score$ חדשות:

$$\begin{aligned} score_2(j) &= -\frac{\langle v^j - \bar{v}^j \cdot 1_m, y - \bar{y} \cdot 1_m \rangle^2}{\|v^j - \bar{v}^j \cdot 1_m\|_2^2 \cdot \|y - 1_m \cdot \bar{y}\|_2^2} \\ &\Downarrow \\ score_3(j) &= \left| \frac{\langle v^j - \bar{v}^j \cdot 1_m, y - \bar{y} \cdot 1_m \rangle}{\|v^j - \bar{v}^j \cdot 1_m\|_2 \cdot \|y - 1_m \cdot \bar{y}\|_2} \right| = \left| \frac{\frac{1}{m} \langle v^j - \bar{v}^j \cdot 1_m, y - \bar{y} \cdot 1_m \rangle}{\sqrt{\frac{1}{m} \|v^j - \bar{v}^j \cdot 1_m\|_2^2} \cdot \sqrt{\frac{1}{m} \|y - 1_m \cdot \bar{y}\|_2^2}} \right| \\ &\Downarrow \\ score_3(j) &= \left| \frac{\frac{1}{m} \sum_{i=1}^m ((x_i)_j - \bar{v}^j) (y_i - \bar{y})}{\sqrt{\frac{1}{m} \sum_{i=1}^m ((x_i)_j - \bar{v}^j)^2} \cdot \sqrt{\frac{1}{m} \sum_{i=1}^m (y_i - \bar{y})^2}} \right| \end{aligned}$$

כאשר הפונקציה האחרונה למעשה מדרגת הפוך מהראשונה - $features$ עם ציון גבוה הם משמעותיים יותר. למעשה מה שקיבלנו הוא פונקציית $Pearson's correlation coefficient$ אשר מודדת קורלציה לינארית בין שני משתנים. ערכיה של פונקציה זו הם בין -1 ל- 1 כאשר 1 מציינ קורלציה חיובית מושלמת ו- -1 קורלציה שלילית מושלמת.

כאשר נרצה ללמד מפרידים לינאריים במימד $feature\ space$ גבוה נראה כי המימד מגדיל הן את ה- $sample\ complexity$ והן את ה- $computational\ complexity$. שיטת $SVM - support\ vector\ machine$ באה להתמודד עם אתגר ה- $sample\ complexity$ על ידי מציאת על-מישור מפריד בעל המרחק ($margin$) הגדול ביותר מהדגימות בסט האימון



הגדרה 18.1 יהי $S = \{(x_i, y_i)\}_{i=1}^m$ סט אימון של m דוגמאות מעל \mathbb{R}^d ותיוגים $y_i \in \{\pm 1\}$. נאמר כי S ניתן להפרדה לינארית ($linearly\ separable$) אם קיים על-מישור (w, b) כך שמתקיים $\forall i \in [m] \text{ sign}(\langle w, x_i \rangle + b) = y_i$. תנאי שקול: $\forall i \in [m] \ y_i (\langle w, x_i \rangle + b) > 0$

נבחין כי כל על-מישור הניתן להפרדה לינארית מקיים את עיקרון ה- ERM ביחס לפונקציית האיבוד $0 - 1$ ועל כן נרצה להיות מסוגלים לבחור חצי מרחב ספציפי מתוך הקבוצה.

Margin and Hard – SVM 18.1

טענה 18.2 בהינתן על-מישור (w, b) כך ש- $\|w\| = 1$ ונקודה x במרחב, אזי המרחב בין הנקודה לעל-מישור הוא $|\langle w, x \rangle + b|$

הוכחה: נשים לב כי המרחק בין הנקודה לעל-מישור מוגדר בתור $\min \{ \|x - v\| : \langle w, v \rangle + b = 0 \}$. על כן אם נבחר $v = x - (\langle w, x \rangle + b)w$ נקבל:

$$\begin{aligned} \langle w, v \rangle + b &= \langle w, (x - (\langle w, x \rangle + b)w) \rangle + b = \langle w, x \rangle - \langle w, (\langle w, x \rangle + b)w \rangle + b \\ &= \langle w, x \rangle - (\langle w, x \rangle + b) \langle w, w \rangle + b = \langle w, x \rangle - (\langle w, x \rangle + b) \|w\|^2 + b \\ &= \langle w, x \rangle - \langle w, x \rangle - b + b = 0 \end{aligned}$$

ובנוסף $\|x - v\| = |\langle w, x \rangle + b| \cdot \|w\| = |\langle w, x \rangle + b|$. כל כן המרחק הינו לכל היותר $|\langle w, x \rangle + b|$. כעת בהינתן נקודה u

על העל-מישור אזי $\langle w, u \rangle + b = 0$ ונקבל:

$$\begin{aligned}
 \|x - u\|^2 &= \|x - v + v - u\|^2 \\
 &= \|x - v\|^2 + \|v - u\|^2 + 2 \langle x - v, v - u \rangle \\
 &\geq \|x - v\|^2 + 2 \langle x - v, v - u \rangle \\
 &= \|x - v\|^2 + 2 \langle x - x + (\langle w, x \rangle + b) w, v - u \rangle \\
 &= \|x - v\|^2 + 2 (\langle w, x \rangle + b) (\langle w, v \rangle - \langle w, u \rangle) \\
 &= \|x - v\|^2
 \end{aligned}$$

כל כן המרחק בין x ל- u הוא לפחות כמו המרחק בין x ל- v אשר ראינו כי הוא: $|\langle w, x \rangle + b|$ כנדרש. ■

הגדרה 18.3 המרחק בין הנקודה הקרובה ביותר ב- S לעל-מישור המפריד $\{u : \langle w, u \rangle + b = 0\}$ נקרא ה- $margin$ ושווה ל:

$$\min_{i \in [m]} |\langle w, x_i \rangle + b|$$

הגדרה 18.4 כלל למידה $Hard - SVM$ הינו:

$$\underset{(w,b) : \|w\|=1}{argmax} \min_{i \in [m]} |\langle w, x_i \rangle + b| \quad s.t. \quad \forall i, y_i (\langle w, x_i \rangle + b) > 0$$

הגדרה 18.5 בעיית אופטימיזציה תקרא $Quadratic Programming$ אם היא ניתנת לכתיבה בצורה הבאה:

$$\min_{w \in \mathbb{R}^n} \left(\frac{1}{2} w^T Q w + c^T w \right) \quad s.t. \quad A w \leq d$$

כאשר $Q \in \mathbb{R}^{n \times n}$, $A \in \mathbb{R}^{m \times n}$, $c \in \mathbb{R}^n$, $d \in \mathbb{R}^m$

כעת, בדומה לפתרון בעיות של תכנון לינארי אשר קיימים אלגוריתמים יעילים לפתרון, נראה כיצד להציג את הבעיה כבעיית תכנון $quadratic$ אשר גם לה קיימים אלגוריתמים יעילים לפתרון.

טענה 18.6 יהיו v^*, c^* פתרונות אופטימליים של:

$$(1) \underset{(w,b)}{argmin} \|w\|^2 \quad s.t. \quad \forall i y_i (\langle w, x_i \rangle + b) \geq 1$$

נסמן $\gamma = \frac{1}{\|v^*\|}$ אזי $w^* = \gamma v^*$, $b^* = \gamma c^*$ פתרונות אופטימליים לכלל הלמידה של $Hard - SVM$ הלינארי.

הוכחה: תחילה נראה כי w^*, b^* הינם אכן פתרונות של שלל הלמידה $Hard - SVM$ (כבעיית תכנון לינארי) ולאחר מכן נוכיח אופטימליות. יהי w פתרון חוקי לכלל הלמידה $Hard - SVM$ אזי w אכן מגדיר על-מישור מפריד ומתקיים כי $|\langle w, x_i \rangle + b| = y_i (\langle w, x_i \rangle + b)$ בנוסף, בהיותו w פתרון חוקי, תנאי החיוביות מתקיים ולכן נוכל לייצג את הכלל בתור:

$$(2) \operatorname{argmax}_{||w||=1, b} \min_i y_i (\langle w, x_i \rangle + b)$$

כעת נתבונן בפתרונות v^*, c^* . בהיותם אופטימליים עבור (1) מתקיים כי:

$$\min_i y_i (\langle v^*, x_i \rangle + c^*) = 1$$

שכן אחרת נוכל לחלק את v^* במספר חיובי ולקבל פתרון חוקי קטן מהמינימלי. נשים לב כי w^*, b^* פתרון חוקי עבור (2) שכן $||w^*|| = ||\gamma v^*|| = \left\| \frac{1}{||v^*||} v^* \right\| = 1$ ומתקיים:

$$\min_i y_i (\langle \gamma v^*, x_i \rangle + \gamma c^*) = \gamma \left(\min_i y_i (\langle v^*, x_i \rangle + c^*) \right) = \gamma$$

(נבחין כי γ איננו תלוי ב- i ולכן נוכל להוציאו מהמכפלה הפנימית, להוציא כגורם משותף ולבסוף להוציא מחוץ ל- \min).

כעת נראה כי זהו פתרון אופטימלי עבור (2). יהיו $w_2 \neq v^*, b_2 \neq c^*$. יהיו המקיימים

$$\min_i y_i (\langle w_2, x_i \rangle + b_2) = \delta > \gamma$$

כעת נוכל להגדיר פתרון $\frac{w_2}{\delta}, \frac{b_2}{\delta}$ אשר חוקי עבור (1) ועבורו מתקיים כי $\frac{1}{\delta} < \frac{1}{\gamma}$ וקיבלנו פתרון טוב יותר מ- v^*, c^* .
 בסתירה לאופטימליות שלהם. ■

18.1.1 המקרה ההומוגני

נוכל לעשות רדוקציה מהמקרה של למידת חצאי-מרחב לא הומוגניים ללמידת חצאי-מרחב הומוגניים. נעשה זאת על ידי הוספת $feature$ נוסף לכל דגימה x_i . כאשר נעשה זאת נעלה את המימד ל- \mathbb{R}^{d+1} . נשים לב כי במקרה הקודם ה- $bias$ b לא נכלל בהגבלה של $||w||^2 = 1$ ואם נלמד את הבעיה ב- \mathbb{R}^{d+1} למעשה אנו מכניסים רגולריזציה לגודל ה- $bias$.

18.1.2 Sample Complexity של $Hard - SVM$

נזכר כי ה- $VCdim$ של חצאי מרחב ב- \mathbb{R}^d הינו $d + 1$. מכאן כי ה- $sample complexity$ של למידת החצאי-מרחב גדלה עם המימד. מהמשפט המרכזי של הלמידה הסטטיסטית אנחנו מקבלים כי אם מספר הדגימות קטן משמעותית מ- $\frac{d}{\epsilon}$ אז לא קיים אלגוריתם שיכול ללמוד את החצי-מרחב בדיוק ϵ .

כדי להתמודד עם הבעיה נצטרך להוסיף הנחה עבור התפלגות המידע. נגדיר הנחה כי המידע ניתן להפרדה לינארית עם $margin$ של γ ונראה שבתנאים אלה ה- $sample complexity$ חסום מלמעלה על ידי $\frac{1}{\gamma^2}$.

הגדרה 18.7 יהי \mathcal{D} התפלגות מעל $\mathbb{R}^d \times \{\pm 1\}$. נאמר כי \mathcal{D} ניתנת להפרדה לינארית עם $(\gamma, \rho) - margin$ אם קיים (w^*, b^*) כך שמתקיים: $||w^*|| = 1, \mathbb{P}_{(x,y) \sim \mathcal{D}} (y (\langle w^*, x \rangle + b^*) \geq \gamma) = 1, ||x|| \leq \rho$.

באופן דומה נגדיר כי \mathcal{D} ניתנת להפרדה לינארית עם $\gamma - \text{margin}$ עם חצי-מרחב הומוגני אם קיים $(w^*, 0)$ כך שהתנאים הקודמים מתקיימים.

משפט 18.8 תהי \mathcal{D} התפלגות מעל $\mathbb{R}^d \times \{\pm 1\}$ הניתנת להפרדה לינארית עם $\gamma - \text{margin}$. מתקיים כי בהסתברות של לפחות $1 - \delta$ עבור סט אימון מגודל m שגיאת ה-1 של Hard-SVM תהיה לכל היותר:

$$\sqrt{\frac{4\left(\frac{\rho}{\gamma}\right)^2}{m}} + \sqrt{\frac{2\log\left(\frac{2}{\delta}\right)}{m}}$$

18.2 Soft-SVM

על מנת להשתמש ב- Hard-SVM עלינו להניח כי המידע ניתן להפרדה לינארית. כדי להחליש הנחה זו נעשה שימוש ב- Soft-SVM . אם כך נגדיר $\xi_1, \dots, \xi_m \geq 0$ אשר:

$$\underbrace{y_i(\langle w, x_i \rangle + b) \geq 1}_{\text{Hard-SVM}} \Rightarrow \underbrace{y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i}_{\text{Soft-SVM}}$$

אם כך ב- Soft-SVM קיים tradeoff בין ה- margin שמחפשים לבין ה- $\text{violation of constraints}$ שמבוטא ב- ξ_i . tradeoff זה מגולם בפרמטר λ .

הגדרה 18.9 כלל למידה Soft-SVM הינו:

$$\min_{w, b, \xi} \left(\lambda \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right) \text{ s.t. } \forall i, y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i, \xi_i \geq 0$$

הגדרה 18.10 בהינתן (w, b) , סט אימון S ופונקציית איבוד hinge ממוצעת $L_S^{\text{hinge}}((w, b))$ (כאשר $\ell^{\text{hinge}}((w, b), (x, y)) = \max\{0, 1 - y(\langle w, x \rangle + b)\}$) נוכל להגדיר

$$\min_{w, b} \left(L_S^{\text{hinge}}((w, b)) + \lambda \|w\|^2 \right)$$

אשר הינה בעיה שקולה לכלל ה- Soft-SVM . נשים לב כי תחת הגדרה זו אנו רואים כי Soft-SVM נופל תחת הפרדיגמה של $\text{Regularized loss minimization (RLM)}$ וכי $\mathcal{R}(w) = \lambda \|w\|^2$.

באופן כללי ה- VCdim של למידת halfspace תלוי במימד d ועל כן ה- sample complexity גדל יחד עם גדילת d . בניגוד לזה, שימוש ב- SVM מגדיר חסם ל- sample complexity על ידי $\left(\frac{\rho}{\gamma}\right)^2$ אשר יעיל מאוד במצבים שבהם $d \gg \left(\frac{\rho}{\gamma}\right)^2$.

18.3 SVM vs. Boosting

	SVM	Boosting
output classifier	$\text{sign}(\mathbf{w} \cdot \psi(\mathbf{x}))$	$\text{sign}(\mathbf{w} \cdot \psi(\mathbf{x}))$
requirement from ψ	$\langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle$ easy to compute	base classifier is WL
coordinates of ψ	\mathbb{R}	$\{-1, 1\}$
norm on weights	$\ \mathbf{w}\ _2$	$\ \mathbf{w}\ _0$

19 Kernels

ראינו אפוא כי חיפוש מפרידים לינאריים בממדים גבוהים מספק יכולת ביטוי רבה לאלגוריתמים לחיפוש חצאי-מרחב. שיטה זו עולה לנו הן ב- $sample\ complexity$ והן ב- $computational\ complexity$. SVM נתן מענה לבעיית ה- $sample\ complexity$ באמצעות $margin$. על מנת להתמודד עם $computational\ complexity$ נשתמש בשיטת $kernels$.

19.1 Embedding into Feature Space

נראה כי מפרידים לינאריים הינם בעלי יכולת פתרון בעיות מוגבלת שכן יתכן ובמימד הנוכחי לא ניתן להפריד את המידע בצורה טובה. על כן נוכל תחילה למפות את הדגימות ל- $feature\ space$ אחר ובו לחפש מפריד לינארי טוב. הפרדיגמה הכללית למיפוי זה:

1. עבור \mathcal{X} domain ומשימת אימון נבחר פונקציית מיפוי $\psi : \mathcal{X} \rightarrow \mathcal{F}$.

2. בהינתן סט אימון $S = \{(x_i, y_i)\}_{i=1}^m$ ניצור $\hat{S} = \{(\psi(x_i), y_i)\}_{i=1}^m$.

3. נאמן h לינארי מעל \hat{S} .

4. עבור דגימת בדיקה x ניתן חיזוי $h(\psi(x))$.

ההצלחה של שיטה זו תלויה בהצלחה בבחירת פונקציית מיפוי כך שבתמונתה הדוגמאות יהיו $linear\ separable$. דוגמה אחת לפונקציית מיפוי גרית הינה $polynomial\ mapping$. אם במפרידים לינאריים דיברנו על $\langle w, x \rangle$ אזי נוכל להרחיב ל- $p(x)$ כאשר p הינו $multivariate\ polynomial$ מדרגה k .

הגדרה 19.1 $multivariate\ polynomial$ מדרגה k מ- \mathbb{R}^n ל- \mathbb{R} הינו ביטוי מהצורה:

$$p(x) = \sum_{J \in [n]^r : r \leq k} \left(w_J \prod_{i=1}^r x_{J_i} \right)$$

19.2 Kernel Trick

השימוש בפונקציית מיפוי זו אשר מעלה את המימד של ה- $data$ טומנת בתוכה גם בעיות: מציאת ψ איננה בהכרח משימה פשוטה, חישוב $\langle w, \psi(x) \rangle$ עלול להיות יקר וה- $sample\ complexity$ שהוא $k+1$ $VCdim(halfspaces\ in\ \mathbb{R}^k)$ מצריך מספר רב של דגימות אימון.

את בעיית ה- $sample\ complexity$ ניתן לפתור כפי שראינו ב- SVM . על מנת להתמודד עם העלות חישוב (computational complexity) נראה את ה- $Kernel\ Trick$.

שיטה זו הינה יעילה כאשר חישוב המיפוי $\psi(x)$ יקר אבל המכפלה הפנימית $\langle \psi(x_i), \psi(x_j) \rangle$ זול.

משפט 19.2 (representers theorem) תהי $\psi : \mathcal{X} \rightarrow \mathcal{F}$ וכלל הלמידה

$$w^* \in \underset{w}{\operatorname{argmin}} f(\langle w, \psi(x_1) \rangle, \dots, \langle w, \psi(x_m) \rangle) + R(\|w\|)$$

כאשר $f : \mathbb{R}^m \rightarrow \mathbb{R}$ ו- $R : \mathbb{R}_+ \rightarrow \mathbb{R}$ פונקציה מונוטונית עולה, אזי קיים פתרון w^* וקטור $\alpha \in \mathbb{R}^m$ כך ש: $w^* = \sum_{i=1}^m \alpha_i \psi(x_i)$

הערה 19.3 כל כללי ה- SVM מהפרק הקודם (ואופטימיזציות נוספות שניתן לראות בספר הקורס) מקיימות את כלל הלמידה הנ"ל. בנוסף R לדוגמה עבור $Soft-SVM$ הינה למעשה λa^2 או עבור $Hard-SVM$ הינה למעשה a^2 .

הוכחה: יהי w^* פתרון אופטימלי לכלל הלמידה. נוכל לכתבו בתור

$$w^* = \sum_{i=1}^m \alpha_i \psi(x_i) + u$$

כאשר $\langle u, \psi(x_i) \rangle = 0$ לכל i . נגדיר $w = w^* - u$. נראה כי מתקיים $\|w\|^2 = \|w^*\|^2 + \|u\|^2$. נראה כי מתקיים $\|w\| \leq \|w^*\|$. ממנוטוניות R מתקיים כי $R(\|w\|) \leq R(\|w^*\|)$. בנוסף

$$\begin{aligned} \forall i \quad \langle w, \psi(x_i) \rangle &= \langle w^* - u, \psi(x_i) \rangle = \langle w^*, \psi(x_i) \rangle - \langle u, \psi(x_i) \rangle \stackrel{(*)}{=} \langle w^*, \psi(x_i) \rangle \\ &\Downarrow \end{aligned}$$

$$f(\langle w, \psi(x_1) \rangle, \dots, \langle w, \psi(x_m) \rangle) = f(\langle w^*, \psi(x_1) \rangle, \dots, \langle w^*, \psi(x_m) \rangle)$$

■

ולמעשה הראנו כי $w = \sum_{i=1}^m \alpha_i \psi(x_i)$ פתרון אופטימלי כנדרש.

על בסיס משפט זה נשים לב כי לכל i מתקיים:

$$\begin{aligned} \langle w, \psi(x_i) \rangle &= \left\langle \sum_j \alpha_j \psi(x_j), \psi(x_i) \right\rangle = \sum_{j=1}^m \alpha_j \langle \psi(x_j), \psi(x_i) \rangle \\ \|w\|^2 &= \left\langle \sum_j \alpha_j \psi(x_j), \sum_j \alpha_j \psi(x_j) \right\rangle = \sum_{i,j=1}^m \alpha_i \alpha_j \langle \psi(x_i), \psi(x_j) \rangle \end{aligned}$$

אזי תהי $K(x, x') = \langle \psi(x), \psi(x') \rangle$ פונקציית *kernel* נוכל לכתוב את כלל הלמידה בתור

$$\min_{a \in \mathbb{R}^m} f \left(\sum_{j=1}^m \alpha_j \langle \psi(x_j), \psi(x_1) \rangle, \dots, \sum_{j=1}^m \alpha_j \langle \psi(x_j), \psi(x_m) \rangle \right) + R(\|w\|)$$

$$\Downarrow$$

$$\min_{a \in \mathbb{R}^m} f \left(\sum_{j=1}^m \alpha_j K(x_j, x_1), \dots, \sum_{j=1}^m \alpha_j K(x_j, x_m) \right) + R \left(\sqrt{\sum_{i,j=1}^m \alpha_i \alpha_j K(x_i, x_j)} \right)$$

וביטוי זה איננו תלוי כלל ב- w וכל שעלינו לדעת זה כיצד לחשב את המטריצה $G \in M^{m \times m}(\mathbb{R})$ כך ש: $G_{i,j} = K(x_i, x_j)$ ולמעשה מתקיים כי:

$$\langle w, \psi(x_i) \rangle = (G\alpha)_i$$

$$\underbrace{\|w\|^2}_{\Downarrow} = \alpha^T G \alpha$$

$$\Downarrow$$

$$\operatorname{argmin}_{a \in \mathbb{R}^m} (f(G\alpha) + \lambda \alpha^T G \alpha)$$

היתרון בעבודה עם *kernels* על פני מציאת w אופטימלי ב-*feature space* הוא שבמקרים שמימד ה-*feature space* מאוד גבוה מימוש וחישוב *kernel* יכול להיות פשוט (וזול חישובית) בהרבה.

דוגמה: כאשר ניקח את כלל הלמידה של *Soft-SVM* ונבטאו באופן זה נקבל:

$$\min_{w,b,\xi} \left(\lambda \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right) \quad s.t \quad \forall i, y_i (\langle w, x_i \rangle + b) \geq 1 - \xi_i \quad \xi_i \geq 0$$

$$\Updownarrow$$

$$\min_{a \in \mathbb{R}^m} \left(\lambda \alpha^T G \alpha + \frac{1}{m} \sum_{i=1}^m \max \{0, 1 - y_i (G\alpha)_i\} \right)$$

משוואה זו נוכל לכתוב בעיית *quadratic programming* עבורה קיים אלגוריתם יעיל. ברגע שנלמד את המקדמים α חישוב דגימה חדשה מתבצע על ידי:

$$\langle w, \psi(x) \rangle = \sum_{j=1}^m \alpha_j \langle \psi(x_j), \psi(x) \rangle = \sum_{j=1}^m \alpha_j K(x_j, x)$$

באופן כללי אפיון של פונקציות *kernel* מתבצע באמצעות המשפט הבא:

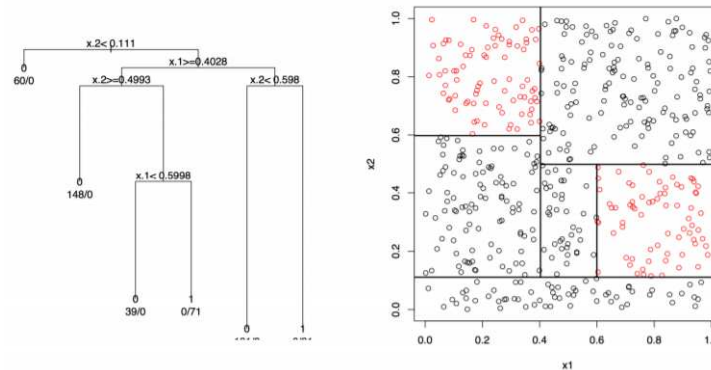
הגדרה 19.4 נאמר כי מטריצה סימטרית $M \in M^{n \times n}(\mathbb{R})$ הינה *positive semidefinite* אם ורק אם לכל וקטור $z \in \mathbb{R}^n$ מתקיים כי $z^T M z \geq 0$.

משפט 19.5 פונקציה סימטרית $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ מממשת מכפלה פנימית במרחב *Hilbert* אם ורק אם היא *semidefinite positive*. כלומר אם לכל x_1, \dots, x_m מטריצת גראם $G_{i,j} = K(x_i, x_j)$ הינה מטריצה *positive semidefinite*.

Decision Trees 20

עץ החלטה הוא predictor $h : \mathcal{X} \rightarrow \mathcal{Y}$ אשר נותן חיזוי לתיגים על ידי מסלול בעץ החלטה מהשורש לעלים. בכל קודקוד ישנו תנאי התלוי ב-*feature* אחד או יותר ובהתאם לקיומו נבחר תת העץ אליו ממשיכים. בדרך הכלל התנאי הינו תנאי של *threshold* (\leq, \geq) עבור *feature* יחיד.

נוכל לחשוב על עץ החלטה בתור חיתוכים של המרחב שלנו כל פעם לחצאי-מרחב בהתאם לתנאי ולבסוף לתאים כאשר כל עלה ממפה לתא ב-*feature space*. נובע כי עץ בעל k עלים יכול לנתח סט $C \subset \mathcal{X}$ בגודל k . מכאן כי עצים בגדלים שונים מספקים מחלקות היפותזות בעלות *VCdim* שונות.



נבחין כי גם במקרה הפשוט של $\mathcal{X} = \{0, 1\}^d$, $\mathcal{Y} = \{0, 1\}$ ותנאי החלטה עם *feature* בודד "מינה אם $1_{x.i=1}$ " אם ניקח עץ עם כל האופציות, 2^d עלים, ועומק $d + 1$ נקבל כי $VCdim(\mathcal{H}) = 2^d$. כדי להתגבר על בעיה זאת, ובכל גם למנוע מצב של *overfitting*, נוכל להשתמש בעיקרון *minimum description length (MDL)* שנמצא בפרק 7 בספר הלימוד.

באופן כללי עבור d *features* ועצים בינאריים נקבל כי בהסתברות לפחות $1 - \delta$ עבור m דגימות, לכל n טבעי המייצג את מספר הקדקודים בעץ, ולכל עץ $h \in \mathcal{H}$ נקבל כי:

$$L_{\mathcal{D}}(h) \leq L_S(h) + \sqrt{\frac{(n+1) \log_2(d+3) + \log\left(\frac{2}{\delta}\right)}{2m}}$$

20.1 Bias – Complexity tradeoff

אם כך נשים לב כי במקרה של עצי החלטה המנופים ב- $tradeoff$ זה הם למעשה מספר העלים (שקול לעומק העץ) שנבחר. עבור מודלים של עצי החלטה עם תנאים מסובכים יותר או עצים שאינם בינאריים הן העומק המקסימלי והן המפתח של כל קודקוד הם המנופים שב- $tradeoff$.

20.2 אלגוריתמים

נרצה אלגוריתם שיחזיר עצים הממזערים את צד ימין בחסם האיבוד שתואר לעיל. אולם היות ומדובר בחיפוש בכל תתי הקבוצות של d כתלות ב- n ו- m לא קיים אלגוריתם יעיל לכך ולכן האלגוריתמים השונים למציאת עצי החלטה מבוססים על יוריסטיקות שונות.

המודל הכללי לבניית העץ מתחיל בעץ בעל עלה אחד (השורש) אשר תיוגו בהתאם ל- $majority vote$ ב- $training set$. כעת נרוץ מספר איטרציות כאשר בכל אחת נבחן את ההשפעה של פיצול עלה ספציפי. נגדיר מדד שיכמת את השיפור שבפיצול. כעת על בסיס כל הפיצולים האפשריים ניקח את הפיצול הממקסם את השיפור או שנבחר שלא לפצל את העלה.

כדי לסווג באמצעות העץ יהיו $\mathbb{R}^d = \bigcup_{j=1}^N R_j$ ה- $partitions$ שהתקבלו על ידי על החלטה שבנינו אזי בהינתן דגימה x לתיוג תיוגה יהיה:

$$\hat{p}_y(R_j) = \frac{1}{n_j} \sum_{x_i \in R_j} \mathbb{1}_{y_i=y}$$

$$\underset{y=\pm 1}{\operatorname{argmax}} \hat{p}_y(R_j) \quad s.t \quad \text{and } n_j = \text{number of training points in } R_j$$

20.2.1 CART – Classification and Regression Trees

באלגוריתם זה היוריסטיקה המנחה את בניית העץ הינה חמדניות. נתחיל מהשורש ונרד לעלים, כאשר בכל שלב נבחר קואורדינטה $1 \leq j \leq d$ ו- $c \in \mathbb{R}$ כך שהפיצול לפי $x_j > c$ יצור את העץ עם ה- $classification error$ הנמוך ביותר. כעת נמשיך לעשות זאת לכל מלבן שהתקבל ונחפש את האופציה שבה השיפור המשמעותי ביותר ב- $misclassification error$. כלומר:

יהי $S = \{(x_i, y_i)\}_{i=1}^m$ סט האימון ונניח כי חילקנו את ה- $feature space$ ל- M אזורים R_1, \dots, R_M . נניח כי התשובה עבור כל אזור היא הקבוע $c_m \in \mathcal{Y}$ אזי

$$f(x) = \sum_{m=1}^M c_m I(x \in R^d)$$

עבור קריטריון המינימום של $\sum (y_i - f(x_i))^2$ *minimum sum of squares* נראה כי \hat{c}_m האופטימלי הוא פשוט ממוצע התיגים באזור:

$$\hat{c}_m = \text{avg}(y_i | x_i \in R^d)$$

כעת בהתאם לעיקרון החמדני לכל $1 \leq j \leq d$ יהי $s \in \mathbb{R}$ ערך ה-*threshold* אזי נגדיר

$$R_1(j, s) = \{x \in \mathbb{R}^d | x_j \leq s\}, R_2(j, s) = \{x \in \mathbb{R}^d | x_j > s\}$$

ונחפש j אשר הפיצול לפיו מביא למינימום את השגיאה בשני תתי העצים:

$$\min_{j,s} \left[\min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2 \right]$$

כאשר ספיציפית למקרה של קלסיפיקציה אנו בעצם מחפשים:

$$\min_{j,s} (n_0 [1 - \hat{p}_{y_0}(R_0)] + n_1 [1 - \hat{p}_{y_1}(R_1)])$$

ועל כן המינימיזציה ניתנת לפתרון על ידי

$$\hat{c}_1 = \text{avg}(y_i | x_i \in R_1(j, s)), \hat{c}_2 = \text{avg}(y_i | x_i \in R_2(j, s))$$

נמשיך באופן רקורסיבי.

הערה 20.1 לשם אינטואיציה ניתן להסתכל על בעיית מציאת על ההחלטה כאנלוג לבעיית חיתוך הבד מקורס אלגוריתמים, כאשר בכל איטרציה של בעיית חיתוך הבד אנו מסתכלים על תמונת אזורים R_1, \dots, R_m שונה עד אשר נחליט להפסיק. בבעיה זו במקום חיתוך אופקי או אנכי לפי ערך כלשהו ישנן d אפשרויות חיתוך לפי ערכים כלשהם.

20.3 Pruning (גיזום)

העצים שחוזרים מהאלגוריתמים לבניית העצים, למרות שבעלי *empirical risk* נמוך, נוטים להיות בעלי *true risk* גבוה. כשם שאופציה אחת להגביל זאת היא על ידי מספר האיטרציות לבניית העץ, נוכל גם לבצע תהליך של *pruning* שבו נקטין את העץ בתקווה ל-*empirical risk* דומה. בדרך כלל תהליך ה-*pruning* הולך *bottom-up* מהעלים לשורש.

נרצה להגדיר קריטריון "עלות" הסיבוכ של העץ (*cost complexity*). יהי T_0 העץ שחזר מתהליך בניית העץ אזי נחפש עץ $T \subset T_0$ אשר נחסיר ממנו חלק מה-*nodes* הפנימיים. נסמן קודקודים טרמינליים (עלים) m בהתייחס לאזור ה- R_m ונגדיר $|T|$ כמספר הקודקודים הטרמינליים ב- T . כעת:

$$N_m = \# \{x_i \in R_m\}, \hat{c}_m = \text{avg}(y_i | x_i \in R_m), Q_m(T) = \frac{1}{N_m} \sum_{x_i \in R_m} (y_i - \hat{c}_m)^2$$

אזי קריטריון העלות סיבוכ של העץ:

$$C_\alpha(T) = \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T|$$

$$\Downarrow$$

$$C_\alpha(T) = \sum_{m=1}^{|T|} [1 - \hat{p}_{c_m}(R_m)] + \alpha |T|$$

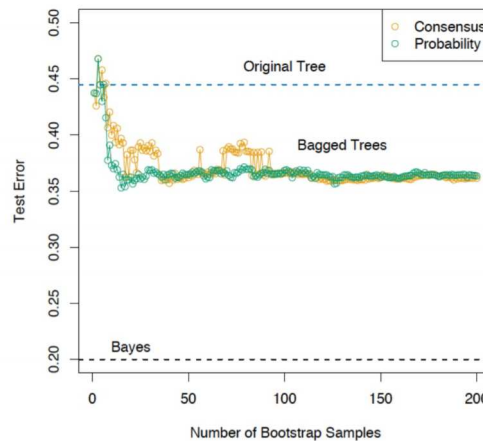
כלומר לכל α נחפש את תת העץ $T_\alpha \subseteq T_0$ אשר מביא למינימום את $C_\alpha(T)$ ומכאן כי α הוא המנוף בין גודל העץ לבין כמה העץ מתאים על ה- $data$. עבור α גדול נקבל עצים קטנים. נראה כי גם אלגוריתם זה נופל תחת פרדיגמת RLM כאשר $\mathcal{R}(T) = \alpha |T|$.

20.3.1 Bias variance

עבור α קטנים מקבל $bias$ נמוך ו- $variance$ גבוה. באופן כללי עצי החלטה הינם בעלי $variance$ גבוה שכן ברגע שבחרנו $feature$ וספ לפצל לפי כל יתר מבנה העץ כבר תלוי בחלטה הזאת.

20.4 שימוש ב- $Bagging$ ו- $Bootstrapping$

כפי שראינו בפרק על $Bagging$ נוכל עבור $B \in \mathbb{N}$ שנקבע לאמן B עצים $\{T_\alpha^{(b)}\}_{b=1}^B$ ובהינתן דגימה חדשה לסווגה על פי ה- $majority vote$ של עצים אלה.



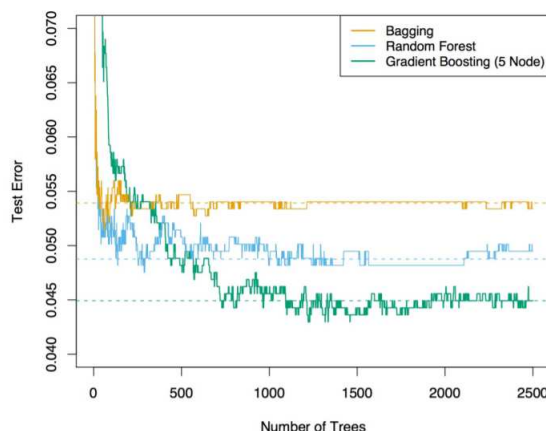
היות ועצי החלטה הינם מודל בעל $variance$ גבוה (החלטת הסף וה- $feature$ בבסיס העץ כבר משנים את מבנה כל העץ) $Bagging$ היא שיטה נוחה להורדת ה- $variance$ מבלי להעלות ב- $bias$ שכן הסיווג מתבצע על סמך מיצוע של עצים שונים.

כדי להבין מדוע זה עובד נסמן p אז סיכוי ה- $misclassification$ של העץ אזי עבור דגימה x ו- B עצים ב- $bagging$ נקבל כי ההסתברות לטעות היא $\mathbb{P}(X \geq \frac{B}{2})$ כאשר $X \sim \text{Bin}(B, p)$. ומכאן שאשר $B \rightarrow \infty$ נקבל כי $\mathbb{P}(X \geq \frac{B}{2}) \rightarrow 0$. נשים לב כי ההנחה הסמויה היא שהעצים בלתי תלויים (המשתנה הבינומי). הנחה זו איננה נכונה שכן העצים אומנו על אותו מדגם ועל כן קיימת קורלציה ביניהם אך עדין נקבל ששימוש ב- $Bagging$ משפר ביצועים כלל ש- B גדל (והיות וקיימת קורלציה לא נוכל באמת להגיע להסתברות 0 לטעות).

20.5 Random Forest

דרך נוספת להוריד את הסיכון לקבלת מסווג עם $true risk$ גבוהה היא על ידי $ensemble$ של עצים קטנים. $random forest$ הוא מסווג המכיל אוסף של $decision trees$. ישנן דרכים רבות ליצור עצים אלה. דרך אחת היא להגדיר $I_1, \dots, I_t \subset [d] \mid |I_i| = k$. כעת ניקח סט אימון S ולכל I_i נדגום את S' אשר הוא m' דגימות בהתפלגות אחידה מעל S עם החזרה. נאמן עץ על סמך ה- $features$ ב- I_i . לבסוף כאשר נרצה לתייג ניקח את התייג על פי ה- $majority vote$ של כלל העצים.

כלל אצבע לבחירת $k: k \sim \sqrt{d}$.



21 k – Nearest Neighbors

בשפחת אלגוריתמים זו, בהינתן דגימה חדשה לסיווג נשתמש ב- $trainset$ כדי למצוא את k הדגימות שהכי קרובות לדגימה שקיבלנו ונשתמש בהן. ההנחה העומדת בבסיס היא שדגימות קרובות מרחבית בעלי תיג דומה.

21.1 האלגוריתם

נסמן $\rho : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ הפונקציה המחזירה את המרחק בין שני אובייקטים ב- \mathcal{X} . יהי $S = \{(x_1, y_1)\}_{i=1}^m \subseteq \mathcal{X}$ ויהי $\pi_1(x), \dots, \pi_m(x)$ סידור הדגימות ב- S על פי מרחקן $\rho(x, x_i)$. כלומר $\rho(x, x_{\pi_i(x)}) \leq \rho(x, x_{\pi_{i+1}(x)})$. $\forall i \in [m]$.

הגדרה 21.1 כלל ה- k NN עבור תיוגים בינאריים מוגדר בתור ה- $majority vote$ של $\{y_{\pi_i(x)} : i \leq k\}$.

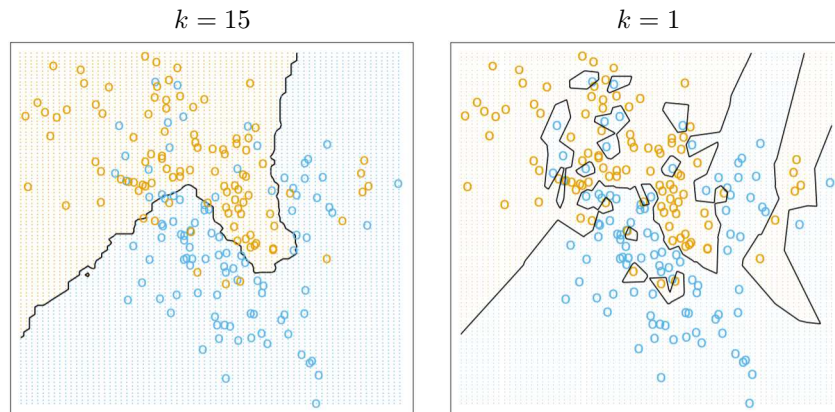
הערה 21.2 במקרה של רגרסיה ($\mathcal{Y} = \mathbb{R}$) נוכל להגדיר את הכלל כממוצע התיוגים של k הדגימות הקרובות: $h_S(x) = \frac{1}{k} \sum_{i=1}^k y_{\pi_i(x)}$

הגדרה 21.3 במקרה הכללי עבור סט תיוגים \mathcal{Y} ופונקציה $\phi : (\mathcal{X} \times \mathcal{Y})^k \rightarrow \mathcal{Y}$ אזי כלל ה- k NN ביחס ל- ϕ הוא:

$$h_S(x) = \phi((x_{\pi_1(x)}, y_{\pi_1(x)}), \dots, (x_{\pi_k(x)}, y_{\pi_k(x)}))$$

21.2 Bias – complexity tradeoff

במודל זה המנוף של $tradeoff$ זה מבוטא בערך ה- k אשר הוא מספר השכנים שממצעים. ככל ש- k קטן יותר אזי המודל "גמיש יותר" ורגיש יותר לרעש שקיים בכל אחת מהדגימות. עבור k גדול יותר המיצוע נעשה על יותר דגימות וכך הרעש של כל דגימה אינדיבידואלית פחות משפיע.



באופן כללי, כל אזור $(neighborhood)$ ב- $NN-k$ תלוי משכנים שמסביב ועל כן הגבולות הללו שבין התיוגים מאוד $wiggly$ ו- $unstable$. מכאן כי סך הכל במודל יש $variance$ גבוהה ו- $bias$ נמוך.

21.3 "Bellman's Curse of dimensionality"

קיימת טענה כי ניתן להגיע ב- $NN-1$ ל- $true\ error$ של $2L_D(h^*) + \epsilon$ כאשר h^* הוא ה- $Bayes\ optimal$ (עמוד 260 בספר הקורס) ובאופן כללי ב- $NN-k$ לחסם של $(1 + \sqrt{8/k})$ פעמים השגיאה של $Bayes\ optimal$. מהוכחת הטענה נובע כי גודל סט האימון הנדרש גדל אקספוננציאלית עם גדילת מימד ה- $features$ וכי עבור התפלגויות שונות זהו תנאי הכרחי למידה:

משפט 21.4 לכל $c > 1$ ולכל כלל למידה L קיימת התפלגות $\{0, 1\} \times [0, 1]^d$ כך שה- $Bayes\ error$ הינה 0 אבל לכל $m \leq (c+1)^d/2$ דגימות אימון ה- $true\ error$ גדול מ- $\frac{1}{4}$.

על כן נרצה להוריד את מימד ה- $feature\ space$ כלל שניתן. ישנן דרכים שונות לעשות זאת כמו שימוש ב- $Kernels$ או $Feature\ Selection$ למציאת ה- $features$ להם המשקל הרב ביותר בניבוי נכון של התיוגים.

22 Convex Learning Problems

22.1 בעיות למידה

להגדרות המתמטיות אודות קבוצות ופונקציות קמורות, לפשיציות וחלקות ראו בפרק ה"דע מוקדם". רוב בעיות הלמידה שנוכל להבטיח עליהן למידה יעילה (PAC וחשובי) הן בעיות קמורות, ליפשיציות וחסומות. ישנן בעיות למידות שאינן קמורות (כמו בשימוש שרשתות נוירונים) ולא כל בעיה קמורה ניתנת ללימוד. בבעיות אלה נתייחס למחלקת ההיפותזות כתת קבוצה $\mathcal{H} \subseteq \mathbb{R}^d$ ונסמן היפותזה במחלקה ב- w .

הגדרה 22.1 בעיית למידה $(\mathcal{H}, \mathcal{Z}, \ell)$ תקרא קמורה אם ורק אם מחלקת ההיפותזות \mathcal{H} הינה קבוצה קמורה ופונקציית האיבוד ℓ קבוצה קמורה בארגומנט הראשון שלה (לכל $z \in \mathcal{Z}$ מתקיים $\ell(w, z)$ היא פונקציה קמורה).

הערה 22.2 על אף שקיימים אלגוריתמים יעילים ללמידה של בעיות קמורות לא מתקיימת הכלה בין הקבוצות: לא כל הבעיות הקמורות למידות ולא כל הבעיות הלמידות קמורות. (קיימת דוגמה של רגרסיה לינארית הומוגנית ספציפית אשר איננה למידה למרות שהיא קמורה - עמוד 164). אולם אם נוסף שני תנאים נוספים נוכל להפוך בעיית למידה זו לקמורה:

1. \mathcal{H} חסומה (עם פרמטר B) - אם לכל $w \in \mathcal{H}$ מתקיים כי $\|w\| \leq B$.

2. ℓ ליפשיצית (עם פרמטר ρ) - אם לכל $z \in \mathcal{Z}$ היא $f(w) = \ell(w, z)$ היא ρ -ליפשיצית. הבעיה תקרא ρ -ליפשיצית.

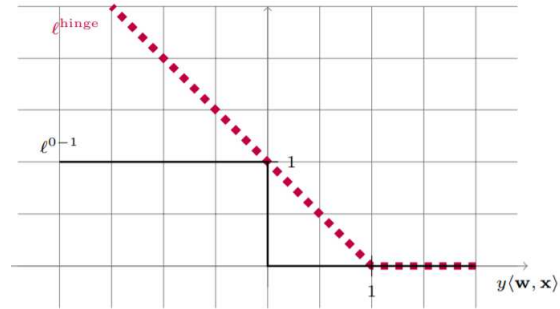
מכאן נובעת ההגדרה הבאה:

הגדרה 22.3 בעיית למידה $(\mathcal{H}, \mathcal{Z}, \ell)$ תקרא קמורה-לפשיצית-חסומה (*Convex – Lipschitz – Bounded*) עם פרמטרים ρ, B אם:

1. מחלקת ההיפותזות \mathcal{H} קבוצה קמורה כאשר $\forall w \in \mathcal{H} \quad \|w\| \leq B$.

2. לכל $z \in \mathcal{Z}$ פונקציית האיבוד $\ell(\cdot, z)$ קמורה ו- ρ -ליפשיצית.

דוגמה: לפונקציית איבוד קמורה הוא ה- ℓ^{hinge} :



הגדרה 22.4 בעיית אופטימיזציה תקרא קמורה אם היא מהצורה $\min_{u \in U} f(u)$ עבור U, f קמורות.

למה 22.5 עבור \mathcal{H} מחלקת היפותזות קמורה ו- ℓ פונקציית איבוד קמורה אזי $ERM_{\mathcal{H}}$, הבאה למינימום של השגיאה האמפירית, הינה בעיית אופטימיזציה קמורה

הוכחה: תהי \mathcal{H} מחלקת היפותזות קמורה ו- ℓ פונקציית איבוד קמורה. נזכר כי $ERM_{\mathcal{H}}(S) = \operatorname{argmin}_{w \in \mathcal{H}} L_S(w)$. עבור סט אימון $S = \{z_i\}_{i=1}^m$ מתקיים כי $L_S(w) = \frac{1}{m} \sum_{i=1}^m \ell(w, z_i)$ היות ו- ℓ קמורה ופונקציות קמורות סגורות תחת סכום וכפל בסקלר אזי גם L_S קמורה. מכאן כי $ERM_{\mathcal{H}}$ בעיית אופטימיזציה קמורה. ■

22.2 פונקציית איבוד Surrogate

כפי שכבר ראינו ניתן ללמוד בעיות למידה קמורות בצורה יעילה ובמקרים אלה נרצה גם פונקציית איבוד קמורות. למשל פונקציית האיבוד ℓ^{0-1} איננה קמורה ומימוש עיקרון ה- ERM עבורה הוא בעיה NP -קשה. נרצה אפוא להגדיר פונקציית איבוד חדשה אשר תספק לנו את התכונות הקמורות ושתחסום מלמעלה את פונקציית האיבוד המקורית. עבור הדוגמה של ℓ^{0-1} פונקציית האיבוד המתאימה לכך היא ℓ^{hinge} .

ברגע שהגדרנו את ה-*surrogate loss function* נרצה לבדוק את חסם השגיאה. עבור הדוגמה של ℓ^{0-1}, ℓ^{hinge} נראה כי:

$$L_{\mathcal{D}}^{hinge}(\mathcal{A}(S)) \leq \min_{w \in \mathcal{H}} L_{\mathcal{D}}^{hinge}(w) + \epsilon \quad \text{for } L_{\mathcal{D}}^{hinge}(w) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell^{hinge}(w, (x, y))]$$

מתכונות ה-*surrogate* נוכל להקטין את אגף שמאל ולכתוב את אגף ימין באופן הבא:

$$L_{\mathcal{D}}^{0-1}(\mathcal{A}(S)) \leq \min_{w \in \mathcal{H}} L_{\mathcal{D}}^{0-1}(w) + \left(\min_{w \in \mathcal{H}} L_{\mathcal{D}}^{hinge}(w) - \min_{w \in \mathcal{H}} L_{\mathcal{D}}^{0-1}(w) \right) + \epsilon$$

וקיבלנו חסם עליון ל- L^{0-1} על ידי שלושה ביטויים:

1. הראשון הוא ה-*approximation error* המודד כמה טוב מצליחה מחלקת ההיפותזות על ההתפלגות.
2. השני הוא ה-*estimation error* (ϵ) הנובע מההבדל בין סט האימון לכלל ההתפלגות.
3. השלישי הוא הביטוי בסוגריים והוא ה-*optimization error*. הוא מודד את ההבדל בין ה-*approximation error* ביחס ל-*surrogate loss* ופונקציית האיבוד המקורית. שגיאה זו נובעת מחוסר היכולת שלנו להביא למינימום את האיבוד ב-*training set* עם פונקציית האיבוד המקורית.

23 Stochastic Gradient Descent

עד כה כאשר עסקנו בבעיות למידה ניסינו להביא למינימום את ה-*empirical risk*. בפרק זה נעסוק בבעיות למידה קמורות ובאמצעות *stochastic gradient descent (SGD)* ננסה להביא למינימום לא את השגיאה האמפירית אלא את ה-*risk function* $L_{\mathcal{D}}(h)$ באופן ישיר. שיטה זו עובדת באופן איטרטיבי כאשר בכל שלב נלך בכיוון הנגדי לגרדיאנט של הפונקציה בנקודה. נשים לב כי היות ואיננו יודעים את \mathcal{D} איננו יודעים את הגרדיאנט. כדי להתגבר על כך *SGD* מתקדם בכיוון שהערך הצפוי בו הוא נגדי לגרדיאנט. היות ובפרק זה נעבוד עם בעיות למידה קמורות נסמן היפותזה במחלקת ההיפותזות בתור $w \in \mathcal{H} \subseteq \mathbb{R}^d$.

23.1 Gradient Descent

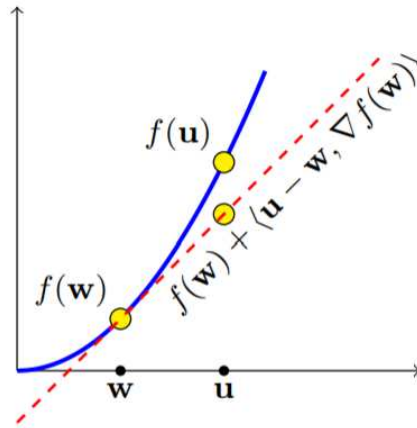
כפי שכבר ראינו הגרדיאנט של פונקציה דיפרנציאבילית $f: \mathbb{R}^d \rightarrow \mathbb{R}$ בנקודה $w \in \mathbb{R}^d$ הינו וקטור הנגזרות החלקיות $\nabla f(x) = \left(\frac{\partial f(x)}{\partial w[1]}, \dots, \frac{\partial f(x)}{\partial w[d]} \right)$. בשלב זה נניח כי הגרדיאנט ידוע לנו

האלגוריתם:

- נתחיל עם $w^{(1)} = 0$
- בכל איטרציה נצעד נגד כיוון הגרדיאנט: $w^{(t+1)} = w^{(t)} - \eta \nabla f(w^{(t)})$ כאשר $\eta > 0$

• פלט האלגוריתם הוא $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

למה 23.1 עבור $f : \mathbb{R}^d \rightarrow \mathbb{R}$ קמורה ואינטגרבילית בנקודה $w \in \mathbb{R}^d$ מתקיים כי $\forall u \in \mathbb{R}^d$ $f(u) \geq f(w) + \langle \nabla f(w), u - w \rangle$. כלומר המשיק של f ב- w נמצא מתחת לגרף הפונקציה.



נראה כי אם w קרוב ל- $w^{(t)}$ אזי הגרדיאנט של f בנקודה w מספק את קירוב טיילור מסדר ראשון של f סביב w על ידי $f(u) \approx f(w) + \langle u - w, \nabla f(w) \rangle$. מכאן כי עבור w הקרוב ל- $w^{(t)}$ נקבל כי

$$f(w) \approx f(w^{(t)}) + \langle w - w^{(t)}, \nabla f(w^{(t)}) \rangle$$

על כן נרצה להביא למינימום את הקירוב וגם לוודא כי אנו נשארים קרובים ל- $w^{(t)}$. פרמטר $\eta > 0$ הוא ששולט ב- $tradeoff$ בין שני אלה ולכן נעדכן את הכלל להיות (אם ניקח כלל זה, נגזור ונשווה לאפס נקבל את הכלל שכתוב באלגוריתם לעיל):

$$w^{(t+1)} = \underset{w}{\operatorname{argmin}} \frac{1}{2} \|w - w^{(t)}\|^2 + \eta \left(f(w^{(t)}) + \langle w - w^{(t)}, \nabla f(w^{(t)}) \rangle \right)$$

23.2 Sub Gradient Descent

אלגוריתם ה- GD מסתמך על כך ש- f דיפרנציאבילית. כעת נקל הנחה זו ונראה כיצד בעזרת $sub - gradient$ של $f(w)$ ב- $w^{(t)}$ נוכל להשתמש באלגוריתם גם על פונקציות שאינן דיפרנציאביליות.

הגדרה 23.2 תהי S קבוצה קמורה פתוחה. פונקציה $f : S \rightarrow \mathbb{R}$ הינה קמורה \iff לכל $w \in S$ קיים v כך ש:

$$\forall u \in S \quad f(u) \geq f(w) + \langle u - w, v \rangle$$

הגדרה 23.3 וקטור v המקיים את התנאי מההגדרה הקודמת יקרא ה- $sub - gradient$ של f בנקודה w . קבוצת כל ה- $sub - gradient$ של f בנקודה w נקראת ה- $differential set$ ומסומנת $\partial f(w)$.

למה 23.4 f הינה קמורה \iff לכל w מתקיים כי $\partial f(w) \neq \emptyset$.

הערה 23.5 בספר יש הסבר על חישוב ה- $sub - gradient$ לפונקציות אולם חלק זה לא הועבר בשיעור. בנוסף האלגוריתם שהוצג עבור שימוש בגרדיאנט מתאים גם עבור שימוש ב- $sub - gradient$.

האלגוריתם

Algorithm 1 Subgradient Descent

Parameter: $T, \eta > 0$
Initialize: $w^{(1)} = \mathbf{0}$
for $t = 1, \dots, T$ **do**
 $w^{(t+1)} = w^{(t)} - \eta v_t$, where $v_t \in \partial f(w^{(t)})$
end for
Return: $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

משפט 23.6 תהי f פונקציה קמורה, ρ -ליפשיצית ויהי $w^* \in \underset{w}{\operatorname{argmin}} f(w)$. אם נריץ את אלגוריתם $sub - gradient descent$ על f ל- T איטרציות עם $\eta = \frac{\|w^*\|}{\rho\sqrt{T}}$ אזי וקטור התוצאה \bar{w} מקיים

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

מסקנה 23.7 תהי f פונקציה קמורה, ρ -ליפשיצית ויהי $w^* \in \underset{w}{\operatorname{argmin}} f(w)$. לכל $\epsilon > 0$ אם נריץ את אלגוריתם $sub - gradient descent$ על f עבור לפחות $\frac{\|w^*\|^2 \rho^2}{\epsilon^2}$ צעדים עם $\eta = \frac{\|w^*\|}{\rho\sqrt{T}}$ אזי וקטור התוצאה \bar{w} מקיים:

$$f(\bar{w}) \leq f(w^*) + \epsilon$$

דוגמה: מציאת על-מישור מפריד. יהי $S = \{(x_i, y_i)\}_{i=1}^m$ אזי נרצה $w \in \mathbb{R}^d$ כך ש: $\forall i \ y_i \langle w, x_i \rangle > 0$. נסמן w^* העל-מישור המפריד של נורמת היחידה ונניח בה"כ כי $\|x_i\| = 1$. נסמן את השולים (margin) של S על ידי $\gamma = \min_i \langle w^*, x_i \rangle$. נחפש את שוליים כך ש: $\max_w \min_i y_i \langle w, x_i \rangle$ השקול למציאת $\min_w \max_i -y_i \langle w, x_i \rangle$ שכן אם w איננו על העל-מישור המפריד ערך המינימום יהיה שלילי ולכן לא יבחר. על כן אנו מתבוננים בבעיה:

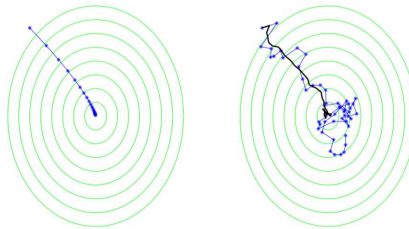
$$\min_w f(w) \text{ where } f(x) = \max_i -y_i \langle w, x_i \rangle$$

נבחין כעת כי f קמורה בהיותה הרכבה של \max על אוסף פונקציות אפיניות (אשר קמורות) וכי ה- $sub - gradient$ של f בנקודה w הוא $-y_i x_i$ עבור i כלשהו השייך ל- $-y_i \langle w, x_i \rangle$.
על כן נרצה את האלגוריתם הבא:

- נאתחל $w^{(1)} = 0$
 - לכל $t = 1, \dots, \left\lceil \frac{1}{\gamma^2} \right\rceil + 1$ נבצע:
 - נמצא i המקיים: $i \in \operatorname{argmax}_i -y_i \langle w, x_i \rangle$
 - נעדכן: $w^{(t+1)} = w^{(t)} + \eta y_i x_i$
 - נחזיר $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$
- נשים לב כי f הינה 1-ליפשיצית וכי $f(w^*) = -\gamma$. מכאן כי לאחר $t > \frac{1}{\gamma^2}$ איטרציות נקבל כי $f(\bar{w}) < f(w^*) + \gamma = 0$ ועל כן \bar{w} הינו אכן על-מישור מפריד כנדרש.
- עלות האלגוריתם - על האלגוריתם לרוץ לפחות $\frac{1}{\gamma^2}$ איטרציות ובכל איטרציה עלות החישוב dm . לכן זמן הריצה של האלגוריתם $O\left(\frac{dm}{\gamma^2}\right)$.

23.3 אלגוריתם Stochastic Gradient Descent (SGD)

נזכר כי אנו מנסים למצוא את ההיפותזה שתביא למינימום את השגיאה (מעל S או מעל \mathcal{D}). SGD הינו אלגוריתם בעל וריאציות רבות אשר המטרה היא שימוש בדגימה אקראית (יחידה או קבוצה, $batch$) לחישוב הגרדיאנט והתקדמות בכיוון הנגדי. בתוחלת נקבל כי אכן נתקרב למינימום. באלגוריתם GD השתמשנו ב- $\nabla f(w^{(t)})$ כאשר f במקרה שלנו היא פונקציית האיבוד. נשים לב כי במקרה של $f(w) = L_{\mathcal{D}}(w) = \mathbb{E}_{z \sim \mathcal{D}} [\ell(w, z)]$ איננו יודעים את $L_{\mathcal{D}}$ ולכן איננו יודעים כיצד לחשב את $\nabla L_{\mathcal{D}}$. בהמשך נראה כיצד להתמודד עם זאת. בנוסף גם אילו ידענו את ההתפלגות והיה ביכולתנו לחשב את הגרדיאנט העלות החישובית של חישובו בכל נקודה הייתה גדולה. על כן ב- SGD , נבחר בכל פעם נקודה $z \sim \mathcal{D}$ (או מספר מצומצם של נקודות) באופן מקרי ונחשב את ה- $sub - gradient$ של $sample$ זה.



23.3.1 SGD עבור מינימום $L_{\mathcal{D}}$

למרות שאיננו יודעים את \mathcal{D} נוכל לחשב שיערוך ל- $sub - gradient$ של $L_{\mathcal{D}}(w)$ היות ומתקיים:

$$\begin{aligned} \nabla L_{\mathcal{D}}(w^{(t)}) &= \nabla \mathbb{E}_{z \sim \mathcal{D}} \ell(w^{(t)}, z) \\ &= \mathbb{E}_{z \sim \mathcal{D}} \nabla \ell(w^{(t)}, z) \end{aligned}$$

וכעת:

Stochastic Gradient Descent (SGD) for minimizing
 $L_{\mathcal{D}}(\mathbf{w})$

parameters: Scalar $\eta > 0$, integer $T > 0$

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

for $t = 1, 2, \dots, T$

sample $z \sim \mathcal{D}$

pick $\mathbf{v}_t \in \partial \ell(\mathbf{w}^{(t)}, z)$

update $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{v}_t$

output $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)}$

מסקנה 23.8 תהי בעיית למידה קמורה-ליפשיצית-חסומה עם פרמטרים ρ, B . אזי לכל $\epsilon > 0$ אם נריץ את SGD עם מספר איטרציות $T \geq \frac{B^2 \rho^2}{\epsilon^2}$ ו- $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$ אזי אלגוריתם SGD מקיים:

$$\mathbb{E}[L_S(\bar{w})] \leq \min_{w \in \mathcal{H}} L_{\mathcal{D}}(w) + \epsilon$$

23.3.2 SGD עבור מינימום L_S (ERM)

נשים לב כי היות ופונקציית הגרדיאנט הינה טרנספורמציה ליניארית מתקיים:

$$\begin{aligned} \nabla L_S(w) &= \nabla \left(\frac{1}{m} \sum_{i=1}^m \ell(w, z_i) \right) \\ &= \frac{1}{m} \sum_{i=1}^m \nabla \ell(w, z_i) \\ &= \mathbb{E}_{i \sim [m]} [\nabla \ell(w, z_i)] \end{aligned}$$

טענה 23.9 אם לכל $i \in [m]$ מתקיים כי $v_i \in \partial \ell(w, z_i)$ אזי $\frac{1}{m} \sum_i v_i \in \partial L_S(w)$. **הוכחה:** נסמן $v = \frac{1}{m} \sum_i v_i$. לכל u מתקיים:

$$\begin{aligned} L_S(u) &= \frac{1}{m} \sum_{i=1}^m \ell(u, z_i) \\ &\geq \frac{1}{m} \sum_{i=1}^m (\ell(u, z_i) + \langle v_i, u - w \rangle) \\ &= L_S(w) + \langle v, u - w \rangle \end{aligned}$$

כלומר אם נגדיל דגימה מסט האימון ונלך נגד כיוון ה- $subgradient$ המתאים לדגימה z_i בנקודה w אזי בתוחלת נתקרב לערך הנכון. ■

אזי אלגוריתם ה- SGD להבאת ה- $L_S(w)$ למינימום נראה כך:

Stochastic Gradient Descent (SGD) for minimizing $L_S(w)$

parameters: Scalar $\eta > 0$, integer $T > 0$

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

for $t = 1, 2, \dots, T$

 choose $i \in [m]$ at random

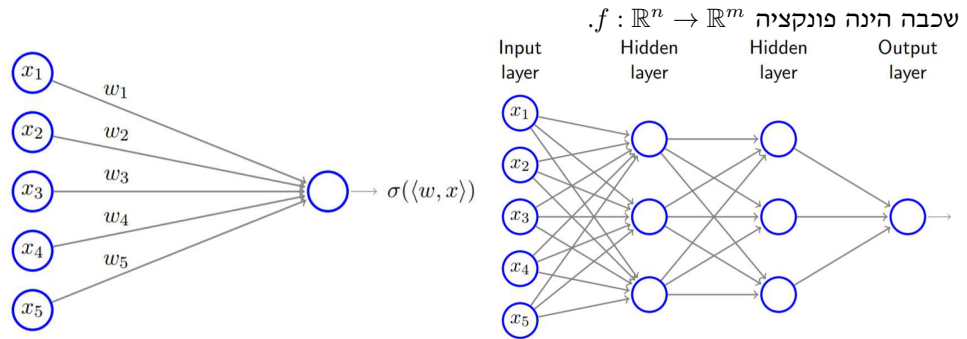
 update $w^{(t+1)} = w^{(t)} - \eta v_t$, where $v_t \in \partial \ell(w^{(t)}, z_i)$

output $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

מסקנה 23.10 תהי בעיית למידה קמורה-ליפשיצית-חסומה עם פרמטרים ρ, B . אזי לכל $\epsilon > 0$ אם נריץ את SGD עם מספר איטרציות $T \geq \frac{B^2 \rho^2}{\epsilon^2}$ ו- $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$ אזי אלגוריתם SGD מקיים:

$$\mathbb{E}[L_S(\bar{w})] - L_S(w^*) \leq \epsilon$$

הרעיון המרכזי ב-*neural networks* הוא הרכבה של צירופים לינאריים של קלט כ-*features* ומידול פונקציית המטרה כפונקציה לא לינארית של *features* אלה. כל אחד מהקודקודים ברשת, הנוירונים, מבצע חישוב פשוט, כאשר ההרכבה שלהם מאפשרת ביצוע חישובים מורכבים. כל רשת מכילה אוסף שכבות כאשר כל שכבה למעשה מבצעת פעולה חישובית כלשהי. סוגי שכבות שונות, ושילובים של פונקציות שונות הם שבסופו של דבר נותנים לרשת הנוירונים את הכוח שלה. כל



בקורס אנו עוסקים רק ברשתות המכונות *feedforward networks*. *feedforward networks* הינן גרף מכוון חסר מעגלים $G = \langle V, E \rangle$ ופונקציית משקולות $w: E \rightarrow \mathbb{R}$. כל *neuron/nodes* מיוצג על ידי פונקציה פשוטה $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ המכונה פונקציית האקטיבציה (*activation function*) של ה-*node*. כל קשת בגרף מחברת בין ה-*output* של *node* קודם ל-*input* של *node* הבא כאשר עבור *node* נתון ה-*input* שלו הוא שיכלול ה-*outputs* של סט *nodes* ומשקולות לכל *node* בסט (כפי שניתן לראות בתמונה השמאלית).

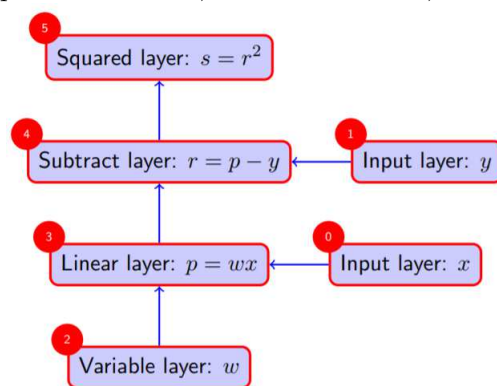
כל ה-*nodes* מאורגנים בשכבות, *layers*, ומכאן כי נוכל לתאר את כל ה-*nodes* בתור: $V = \dot{\bigcup}_{t=0}^T V_t$ וכל קשת ב- E מחברת בין *nodes* ב- V_{t-1} ל- V_t (ניתן להגדיר רשתות אחרות בהן קשתות יכולות להיות תוך שכבתיות). השכבה V_0 מכונה ה-*input layer* והשכבה V_T ה-*output layer*. יתר השכבות V_1, \dots, V_{T-1} מכונות *hidden layers*. בתמונות הנ"ל כל *node* מייצג נוירון אשר מקבל קלט מ- n *input nodes* ומוציא *output* כלשהו על בסיס פונקציית האקטיבציה שמחשב הנוירון. כמובן שנוכל לתת לכל קלט בסט הקלטים של נוירון משקולת w ובכך לציין איזה *input* משמעותי יותר לנוירון (כאשר מדברים על אימון/בניית הרשת מדברים על מהן הפונקציות שמבצעים הנוירונים וכן מהן המשקולות השונות). שתי השכבות החיצוניות של רשת הנוירונים מכונות ה-*input layer* וה-*output layer* אליהן יש לנו גישה. אוסף השכבות הפנימיות מכונות *hidden layers*.

אם כך ניתן לייצג רשת שלמה בתור (V, E, σ, w) אשר מגדירה לנו פונקציה $h_{V, E, \sigma, w}: \mathbb{R}^{|V_0|} \rightarrow \mathbb{R}^{|V_T|}$ כאשר כל סט של פונקציות שכאלה מגדיר מחלקת היפותזות. השלשה (V, E, σ) מכונה הארכיטקטורה של הרשת כאשר בעיות שונות מאופיינות בארכיטקטורה שונה שנמצאה מתאימה להן.

24.1 Computation Graph and Backpropagating

שיטה שעובדת טוב היא שילוב בין SGD ורשתות נוירונים כאשר פונקציית הפלט מהרשת תהיה פונקציית $loss$ שעליה נפעיל SGD להבאת האיבוד למינימום. על כן תהי מחלקת היפותזות המאופיינת בוקטור $\theta \in \mathbb{R}^d$ ופונקציית האיבוד $\ell(\theta, (x, y))$. נניח כי ℓ פונקציה דיפרנציאבילית ב- θ אזי נוכל לחשב את $\nabla \ell(\theta, (x, y))$. כאשר נחפש מינימום של L_D או L_S בעזרת SGD נתחיל מ- $\theta^{(0)}$ ונעדכן בכל איטרציה $\theta^{(t+1)} = \theta^{(t)} - \eta t \nabla \ell(\theta^{(t)}, (x, y))$. בנקודה זו חשוב לשים לב כי בפועל לא מובטח לנו כי הפונקציה או בעיית הלמידה קמורות ובנוסף לא מובטח לנו שקצב הצעדים (η) קטן דיו. על כן יתכן שנתכנס למינימום שאיננו כלל המינימום הגלובאלי (בפונקציה קמורה מינימום מקומי הוא מינימום גלובאלי אך לא כך בפונקציות לא קמורות). בנוסף כלל לא מובטח לנו שהפונקציה הינה דיפרנציאבילית בכלל. בכל זאת, נניח כי הפונקציה ניתנת לתיאור על ידי הרכבה של פונקציות דיפרנציאביליות פשוטות וחיפוש המינימום יעשה באמצעות הגרדיאנט. במבחן התוצאה השיטה עובדת.

אם כך נרצה מודל שיאפשר לנו לחשב את פונקציית ה- $loss$ שנבחר בצורה יעילה וכן את הגרדיאנט בצורה יעילה. היות ואנו מניחים כי הפונקציה הכללית ניתנת לתיאור על ידי הרכבה של פונקציות דיפרנציאביליות פשוטות נוכל על בסיסן, מכלל השרשרת, לחשב את הגרדיאנט של הפונקציה המקורית. כדוגמה נניח כי פונקציית האיבוד ℓ היא $Least Squares$ מעל \mathbb{R} , אזי הגרף מתחת מראה את גרף חישוב של $Least Squares$ במימד אחד.



נראה כי קיימים שלושה סוגים של $layers$:

- $input layer$ המגדירה קבלת קלט כלשהו.
- $variable layer$ המגדירה את הפרמטרים של הבעיה.
- $functional layers$ המגדירות את הפעולות החישוביות הפשוטות שנעשה בכל שלב.

מספורי השכבות מייצגים את הסידור הטופולוגי של השכבות ואת סדר החישוב. ונבחין כי:

1. Forward - נוכל לייצג גרף זה גם בכתוב מתמטי כהרכבה של הפונקציות השונות: $\ell(w) = s(r_y(p_x(w)))$. אם נלך עם הסידור הטופולוגי אכן בסופו נקבל חישוב מלא של ℓ כהרכבה של פונקציות פשוטות. בקוד נוכל לעשות זאת בפשטות על ידי:

For $t = 0, \dots, T - 1$:

$Layer[t].output = Layer[t].function(Layer[t].inputs)$

2. Backward - כעת נרצה לחשב את הנגזרת של פונקציית ℓ . נזכר שעל פי כלל השרשרת $\ell'(w) = s'(r_y(p_x(w))) \cdot r'_y(p_x(w)) \cdot p'_x(w)$ ועל כן אם לאחר שעלינו בעץ עם כיוון הסידור הטופולוגי וחישובנו את ערכי הביניים, נרד בעץ בחזרה ובכל שלב נחשב את הנגזרת של הפונקציה הפשוטה באותה שכבה נצליח לחשב את הנגזרת של פונקציות מסובכות באופן יעיל. זהו ה-backpropagation. בקוד נוכל לעשות זאת בפשטות על ידי:

```

Layer[T - 1].delta = 1
For t = T - 1, ..., 0 :
  For i in Layer[t].inputs :
    i.delta = Layer[t].delta * Layer[t].derivative(i, Layer[t].inputs)

```

δ מייצג את תוצאת חישוב הנגזרת באותו השלב. למעשה אם נבטיח שכשנבנה את גרף החישוב נשתמש רק בפונקציות בעלות נגזרות פשוטות שנוכל לשמור כמשתנים של השכבה וכך בדרך חזרה למטה נצליח לחשב את הנגזרת ביעילות. כשנסיים את ההרצה בגרף למעשה קיבלנו את הערך של $\nabla \ell$, עבור הנקודה w שהזנו כמשתנה. נעשה זאת בכל איטרציה של ה- SGD לחישוב ה- $w^{(t)}$.

דוגמאות ל- $layers$ נפוצים:

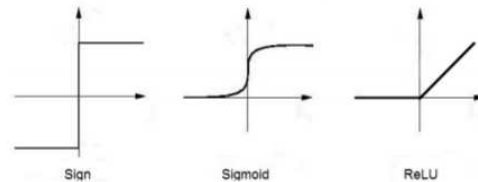
1. שכבה אפינית - $O = WX + b \cdot 1^T$ עבור $W \in \mathbb{R}^{m \times n}$, $x \in \mathbb{R}^{n \times c}$, $\in \mathbb{R}^m$

2. שכבה אונרית - $f : \mathbb{R} \rightarrow \mathbb{R}$, $\forall i, o_i = f(x_i)$, כאשר f לדוגמה היא:

(א) פונקציית סיגמואיד - $f(x) = (1 + \exp(-x))^{-1}$

(ב) פונקציית $ReLU$ (Rectified Linear Unit) - $f(x) = \max\{0, x\}$

(ג) $sign$



3. שכבה בינארית - $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, $\forall i, o_i = f(x_i, y_i)$, כאשר f לדוגמה היא:

(א) חיבור - $f(x, y) = x + y$

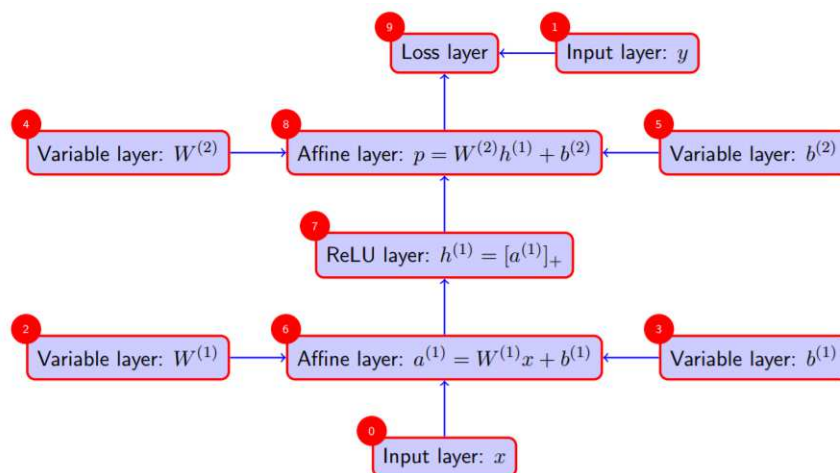
(ב) $hinge$ loss - $f(x, y) = [1 - y_i x_i]_+$

(ג) $logistic$ loss - $f(x, y) = \log(1 + \exp(-y_i x_i))$

עד כה עסקנו ב-Backpropagation של $f : \mathbb{R}^n \rightarrow \mathbb{R}$. נוכל להכליל זאת לפונקציות $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ על ידי כך ש- δ מה-pseudocode יהיה וקטור כך שכמות הקואורדינטות בו שווה לכמות ב- V_T (שכבת ה-output). בנוסף $derivative$ איננו

נגזרת אחת אלא מטריצת היעקוביאן $J_x(f)$ כאשר במיקום i, j -ה נמצאת הנגזרת החלקית של $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ לפי המשתנה j -ה והכפל הינו כפל מטריצות. בכך האלגוריתם למעשה מקיים את *multivariate chain rule*:
 $J_w(f \circ g) = J_{g(w)}(f) J_w(g)$ (במיקום זה ה-*subscript* איננו מייצגת את הקבועים כמו בכתוב של הגרף אלא את המשתנה שלפיו מתבצעת הגזירה).

דוגמה נוספת עבור *computation graph* מורכב יותר הינו הגרף הבא. בגרף זה מבצעים רגרסיה לינארית (השכבה האפיינית התחתונה), מעבירים בשכבת *ReLU* המאפסת ערכים שליליים ואז מבצעים רגרסיה לינארית נוספת.



דוגמה-מנחה - נתבונן בבעיית הלמידה של זיהוי ספרות אזי לתמונות שחור לבן בגודל 28×28 פיקסלים:

$$\mathcal{X} = \{0, \dots, 255\}^{28 \times 28}, \quad \mathcal{Y} = \{0, \dots, 9\}$$

כלומר בעיה זו היא של *multiclassification*. ניקח היפותזה $h : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{Y}|}$ כאשר $h(x)$ הינו וקטור ציונים לכל התיוגים. למשל אם נשתמש ברגרסיה לוגיסטית לשם כך נוכל לייצג זאת כהסתברות לכל אחד מהתיוגים. התיוג שנבחר הינו $\underset{i}{\operatorname{argmax}} h_i(x)$.

- ארכיטקטורת רשת שטובה לבעיה - $x \rightarrow \text{Affine}(500) \rightarrow \text{ReLU} \rightarrow \text{Affine}(10)$ - כלומר עבור תמונה $x \in \mathcal{X}$ נעשה רגרסיה לינארית (עם למשל 500 *nodes* ב-*hidden layer* ראשונה), לאחר מכן נעביר בשכבת *ReLU* להשארת רק ערכים $\mathbb{R}_{\geq 0}$ ואז נבצע רגרסיה שנייה עם 10 *nodes* - אחד לכל ספרה/תיוג.

- שימוש ב-*logistic loss*:

$$\forall i \quad p_i = \frac{\exp(h_i(x))}{\sum_j \exp(h_j(x))} \quad \text{- SoftMax}$$

$$\text{- LogLoss} \quad \text{- עבור } y \text{ התיוג הנכון} \quad \log(p_y) = \log\left(\sum_j \exp(h_j(x) - h_i(x))\right)$$

24.2 Expressiveness and Sample Complexity

נראה כי כוח הביטוי של רשתות נוירונים הוא רב. נוכל לבטא בעזרת רשתות נוירונים קבוצות של "כל הפונקציות מ-A ל-B". נשים לב אבל שביטוי שכזה משמעותו גודל אקספוננציאלי של הרשת כתלות בגודל הקלט.

טענה 24.1 כל פונקציה בוליאנית $f : \{\pm 1\}^n \rightarrow \{\pm 1\}$ ניתנת לביטוי על ידי רשת נוירונים בעלת שכבה $hidden$ יחידה.

משפט 24.2 לכל $n \in \mathbb{N}$ יהי $s(n)$ המספר הטבעי המינימלי כך שקיימת רשת נוירונים בעלת שכבה $hidden$ יחידה עם $s(n)$ נוירונים אשר מממשת את כל הפונקציות מ- $\{0, 1\}^n$ ל- $\{0, 1\}$. אזי $s(n)$ אקספוננציאלי ב- n . **הוכחה:** יהי $G = \langle V, E \rangle$ הגרף המתאים לרשת הנוירונים $\mathcal{H}_{V,E,sign}$ המממשת את כל הפונקציות מ- $\{0, 1\}^n$ ל- $\{0, 1\}$. נראה כי בהינתן $m = 2^n$ וקטורים מעל $\{0, 1\}^n$ נוכל לנתח סט זה. מכאן כי $VCDim(\mathcal{H}_{V,E,sign}) \geq 2^n$. מן הצד השני $VCDim(\mathcal{H}_{V,E,sign})$ חסום על ידי $O(|V| \log |E|) \leq O(|V|^3)$ (משפט בספר בעמוד 274). על כן $|V| \geq \Omega(2^{n/3})$ ועל כן מספר הנוירונים בשכבה ה- $hidden$ אכן אקספוננציאלי ב- n כנדרש. ■

משפט 24.3 יהי $\epsilon \in (0, 1)$ ותהי σ פונקציה סיגמואידית כלשהי. לכל $n \in \mathbb{N}$ יהי $s(n)$ המספר הטבעי המינימלי כך שקיימת רשת נוירונים בעלת $|V| = s(n)$ אשר מחלקת ההיפותזות $\mathcal{H}_{V,E,\sigma}$ יכולה לקרב בדיוק ϵ כל פונקציה 1-ליפשיצית $f : [-1, 1]^n \rightarrow [-1, 1]$. אזי $s(n)$ אקספוננציאלי ב- n .

נראה כי פונקציות הניתנות לחישוב בזמן פולינומיאלי ניתנות לביטוי על ידי רשת בגודל פולינומיאלי:

משפט 24.4 תהי $T : \mathbb{N} \rightarrow \mathbb{N}$ ולכל $n \in \mathbb{N}$ תהי \mathcal{F}_n סט הפונקציות הניתנות לחישוב על ידי מכונת טיורינג בזמן ריצה $T(n)$ אזי קיימים קבועים $b, c \in \mathbb{R}_+$ כך שלכל n קיים גרף מגודל לכל היותר $cT(n)^2 + b$ כך שהרשת $\mathcal{H}_{V_n, E_n, sign}$ מממשת את \mathcal{F}_n .

24.3 Tricks

1. נרמול ה-data - למשל עבור הדוגמה המנחה בתחילת הפרק לחלק את כל הערכים ב-255.
2. נקודת התחלה - היות והבעיה איננה קמורה נקודת ההתחלה חשובה (כאשר קמור אזי מינימום מקומי הינו גלובלי, בבעיה לא קמורה דבר לא מובטח) פתרון שעובד טוב בפועל הוא לאתחל את שורות W (מטריצת המשקלים של השכבה האפניית) באופן רנדומלי בטווח $\left[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right]$.
3. Mini-batches - בכל שלב של SGD נמצע את האיבוד על סמך $k > 1$ דוגמאות רנדומליות. בכך נוריד את ה- $variance$ בעדכון כיוון ההתקדמות ונוכל לממש חלק זה ב- $parallel$ ועל כן לא יעלה הרבה זמן ריצה.
4. Learning rate - נתחיל בערך η כלשהו ונרד כל פעם ב- $\frac{1}{2}$ כאשר האימון מפסיק להתקדם.
5. אלגוריתם SGD - קיימים וריאנטים רבים של האלגוריתם למטרות שונות. נבחר בהתאם למה שנחוץ/טוב לבעיה.

דוגמה ל-*pseudocode* של *SGD* עבור רשתות נוירונים (על *backwards propagation*):

SGD for Neural Networks	Backpropagation
parameters: number of iterations τ step size sequence $\eta_1, \eta_2, \dots, \eta_\tau$ regularization parameter $\lambda > 0$ input: layered graph (V, E) differentiable activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ initialize: choose $\mathbf{w}^{(1)} \in \mathbb{R}^{ E }$ at random (from a distribution s.t. $\mathbf{w}^{(1)}$ is close enough to $\mathbf{0}$) for $i = 1, 2, \dots, \tau$ sample $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$ calculate gradient $\mathbf{v}_i = \text{backpropagation}(\mathbf{x}, \mathbf{y}, \mathbf{w}, (V, E), \sigma)$ update $\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} - \eta_i(\mathbf{v}_i + \lambda \mathbf{w}^{(i)})$ output: $\bar{\mathbf{w}}$ is the best performing $\mathbf{w}^{(i)}$ on a validation set	input: example (\mathbf{x}, \mathbf{y}) , weight vector \mathbf{w} , layered graph (V, E) , activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ initialize: denote layers of the graph V_0, \dots, V_T where $V_t = \{v_{t,1}, \dots, v_{t,k_t}\}$ define $W_{t,i,j}$ as the weight of $(v_{t,j}, v_{t+1,i})$ (where we set $W_{t,i,j} = 0$ if $(v_{t,j}, v_{t+1,i}) \notin E$) forward: set $\mathbf{o}_0 = \mathbf{x}$ for $t = 1, \dots, T$ for $i = 1, \dots, k_t$ set $a_{t,i} = \sum_{j=1}^{k_{t-1}} W_{t-1,i,j} o_{t-1,j}$ set $o_{t,i} = \sigma(a_{t,i})$ backward: set $\delta_T = \mathbf{o}_T - \mathbf{y}$ for $t = T-1, T-2, \dots, 1$ for $i = 1, \dots, k_t$ $\delta_{t,i} = \sum_{j=1}^{k_{t+1}} W_{t,j,i} \delta_{t+1,j} \sigma'(a_{t+1,j})$ output: foreach edge $(v_{t-1,j}, v_{t,i}) \in E$ set the partial derivative to $\delta_{t,i} \sigma'(a_{t,i}) o_{t-1,j}$

Convolutional Networks 24.4

לא נלמד בשלב זה

25.1 הקדמה

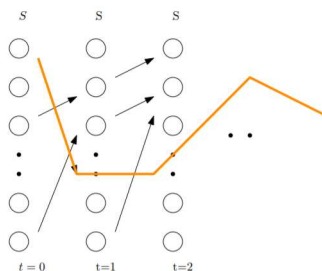
reinforcement learning הינה שיטה חישובית להבנה ואוטומטיזציה של למידה מבוססת *goal* ו-*decision making*. בשיטה זו קיימים מצבים (*situations*) הממפים לפעולות שונות (*actions*), כאשר לבחירת פעולה מסוימת יש *reward/cost*. ה-*agent* (מקביל ל-*learner* בטרמינולוגיה של *supervised learning*) לא מקבל כידע מוקדם איזה פעולות כדאי לו לבחור ועליו ללמוד איזה פעולות יניבו את ה-*reward* המקסימלי. למידה זו שונה מ-*supervised learning*.

ב-*supervised learning* לומדים על בסיס סט דגימות מתווייג על ידי *supervisor* כלשהו ועל הלומד ללמוד על פי איזה *features* וערכי *features* מתקבל התיג הנכון. ב-*reinforcement learning* על ה-*agent* להצליח ללמוד איזה פעולות לבצע (אשר יניבו *reward* מקסימלי) על מצבים שאינם בסט האימון. כלומר על ה-*agent* ללמוד מ"ניסיונו האישי" על הבעיה תוך כדי אינטראקציה עם הסביבה. מכל צורות הלמידה, *reinforcement learning* הינה הדומה ביותר לתהליכי הלמידה שקורים ביצורים חיים ורבים מהאלגוריתמים המרכזיים בשיטה זו הינם בהשראת מערכות למידה ביולוגיות.

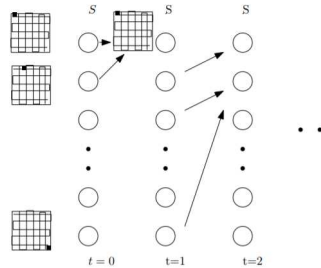
בנוסף לסביבת הבעיה (*environment*) וה-*agent* קיימים מספר אלמנטים מרכזיים ב-*reinforcement learning*:

- פוליסה (*policy*) - ציינו כי בלמידה זו אנו לומדים מיפוי אופטימלי בין מצבים שונים לפעולות שונות. מכאן כי *policy* מגדירה ל-*agent* כיצד להתנהג בכל רגע נתון. למעשה פונקציה $\pi : S \rightarrow A$ כאשר S מתאר את המצבים ו- A את הפעולות. נבחין כי פונקציה זו איננה חייבת להיות דטרמיניסטית ונוכל להגדירה גם בתור $\pi : S \rightarrow 2^A$ כאשר עבור מצב נתון $s \in S$ קיימת הסתברות שונה להחזרת $a \in A$ שונים.
- חיזוקים (*reward signal*) - מגדירים את המטרה (*goal*) שבלמידה. בכל שלב ה-*agent* מקבל מהסביבה *reward* בהתאם למצבים/פעולות שה-*agent* נמצא בהם. מטרתו אם כך למקסם את ה-*reward* בטווח הארוך. זוהי הדרך שלנו לומר ל-*agent* מה אנחנו רוצים להשיג (ולא איך, את זה הוא מגלה לבד). על כן ה-*reward* מאפשר ל-*agent* ללמוד איזה פעולות היו טובות ואיזה לא. כמו בפוליסה, גם כאן ה-*reward* יכול להיות סטוכסטי.
- סונקציית ה-*value* (*value function*) - מגדירה מה למעשה טוב לבצע בטווח הארוך. באופן גס, מייצגת את כמות הרווח הכוללת שה-*agent* יכול לקבל בעתיד. באופן כללי נחפש לבצע *actions* אשר יניבו את ה-*value* המקסימלי ולא בהכרח ה-*reward* המקסימלי.

כלומר, באופן כללי, ב-*reinforcement learning* בכל איטרציה ה-*agent* עושה פעולה, הסביבה משתנה בהתאם וה-*agent* מקבל איזשהו *reward* בהתאם:



דוגמה-מנחה: נדמיין הליכה במבוך כאשר אנו מנסים למצוא מסלול שיאפשר הגעה לצד השני. בכל שלב נתקדם במבוך ונפעל בהתאם. בדוגמה זו מתקיים כי ה-*State space* הם משבצות המבוך וה-*Action space* - הכיוונים שיכול ה-*agent* ללכת. כאשר נתקדם במבוך נראה זאת בתור:



כדי להבין את האופי האיטרטיבי שבו למעשה לא "סיימנו להתאמן" נניח כי מצאנו מסלול אשר מוציא אותנו מהמבוך אבל האם ישנם מסלולים אחרים שמובילים לכך והאם מסלולים אלה מעניקים *reward* גבוה יותר?

ב-*reinforcement learning* המטרה (*goal*) מגולמת בקבלת ה-*reward* מספרי כלשהו בצעדים השונים וה-*goal* הוא מיקסום הסכום של ה-*expected rewards*. נוכל לחשוב על שיטות שונות לעשות זאת:

- Total award - סכימת *reward* כללית $\mathbb{E} \left[\sum_{t=0}^{\infty} r_t \right]$. הבעיה בשיטה זו היא כי יתכן והסכום הינו אינסופי.

- Finite horizon - סכימת *reward* במסגרת $H \in \mathbb{N}$ הצעדים הראשונים: $\mathbb{E} \left[\sum_{t=0}^H r_t \right]$.

- Discounted reward - סכימת *reward* עם משקולת תלויה במספר הצעד $\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$ עבור $\gamma \in [0, 1]$ המכונה גם ה-*discount rate* (נשים לב שעבור $\gamma = 1$ נקבל את *Total award*). שיטה זו נוחה שכן בשונה מ-*total award* פוליסה זו מוגדרת היטב וכן היא מובילה להעדפה של *immdediate reward*.

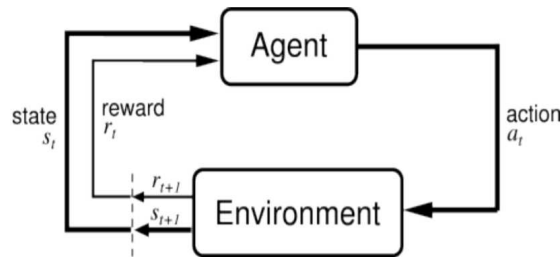
כאשר בונים *agents* ישנן מספר דילמות שונות: האם אנחנו מעוניינים בתוצאות קצרות טווח (*local reward*) או ארוכות טווח (*future reward*); האם אנחנו רוצים שהתהליך יגלה מסלולים חדשים (*exploration*) או יחזיר את התוצאה הכי טובה על מה שכבר מכיר (*exploitation*); ובנוסף איך נדע לקשר בין רווח לבין פעולה שהתרחשה בעבר אשר כעת זיכתה אותנו באותו הרווח (*delay*).

Supervised vs. Reinforcement

Reinforcement learning	Supervised learning	
משוב חלקי (train – and – error)	תשובה נכונה (ground truth)	Feedback מהסביבה
מסלולים (קבוצת המצבים, הפעולות וה-policy)	סט samples בלתי תלויים	קלט ל-learner
החזרת policy. מיפוי בין מצבים ופעולות	החזרת היפותזה. מיפוי בין דוגמאות ותיוגים	פלט
min total cost	min loss	מטרה
בחירת פעולות תביא למצבים שונים	אין	השפעות ארוכות טווח (כלומר האם הדוגמאות הבאות יושפעו מהבחירות שבמצעים)

Finite Markov Decision Processes (MDP) 25.2

Markov decision processes הינם פורמליזציה של תהליך סידרתי של קבלת החלטות, כאשר החלטה שהתקבלה בשלב מסוים עשויה להשפיע על החלטות בנקודות זמן מרוחקות יותר. במקרה של *reinforcement learning* ה-agent מתקשר עם הסביבה באיטרציות $t = 0, 1, 2, 3, \dots$ כאשר בכל איטרציה t ה-agent מחזיר מצב של הסביבה $S_t \in \mathcal{S}$ אשר על בסיס מצב זה בוחר לבצע פעולה $A_t \in \mathcal{A}$. בהתאם לכך מקבל ה-agent את ה-reward $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$ ומוצא עצמו במצב חדש $S_{t+1} \in \mathcal{S}$. כלומר מסלול *MDP (trajectory)* של ה-agent הוא סדרה סופית: $S_0, A_0, R_1, S_1, A_1, \dots$.



באופן פורמלי ההגדרה היא:

הגדרה 25.1 MDP הינו רביעיה $\langle \mathcal{S}, \mathcal{A}, \tau, r \rangle$ כאשר:

1. \mathcal{S} קבוצת מצבים סופית
2. \mathcal{A} קבוצת פעולות סופית
3. $\tau : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]^{|\mathcal{S}|}$ פונקציית הסתברות המעברים. למעשה $\tau(s, a)$ הינה פונקציית המעברים שנותנת הסתברות לכל מצב $s' \in \mathcal{S}$ אם נמצאים במצב $s \in \mathcal{S}$ ומבצעים פעולה $a \in \mathcal{A}$.
4. $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ פונקציית ה-reward.

The Markovian Assumption 25.2.1

עבור MDP סופי נתייחס ל- R_t, S_t כמשתנים מקריים בדידים בעלי הסתברות לקבלם בהתאם ל- S_{t-1} ול- A_{t-1} ומכאן:

הגדרה 25.2 (The Markovian Assumption) ההסתברות בהינתן MDP סופי לקבלת המצב וה- $reward$ הבאים תלוי רק במצב הנוכחי ולא המצבים המוקדמים שהובילו למצב הנוכחי. כלומר:

State Transition

$$\mathbb{P}(S_{t+1} = s' | S_t = s \wedge A_t = a \wedge R_t = r_t \wedge S_{t-1} = s_{t-1} \wedge \dots) = \mathbb{P}(S_{t+1} = s' | S_t = s, A_t = a)$$

Expected Reward

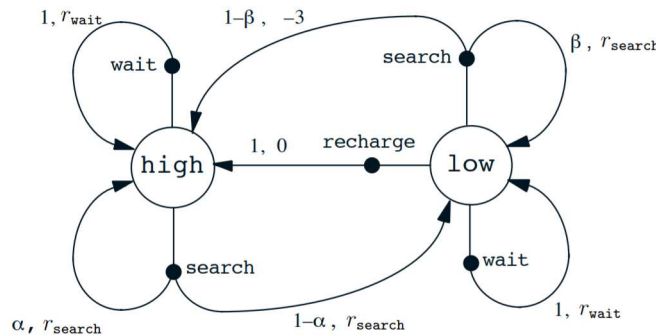
$$\mathbb{E}[R_{t+1} | S_t = s \wedge A_t = a \wedge R_t = r_t \wedge S_{t-1} = s_{t-1} \wedge \dots] = \mathbb{E}[R_{t+1} | S_t = s \wedge A_t = a]$$

נסמן זאת ב:

$$p(s'|s, a) = \mathbb{P}(S_{t+1} = s' | S_t = s, A_t = a)$$

$$r(s, a) = \mathbb{E}[R_{t+1} | S_t = s \wedge A_t = a]$$

נשים לב כי במקרים רבים הנחה זו הינה ראלית. באופן תאורתי נוכל להגדיר מצב כתמונת כל ההיסטוריה כולה (אשר עלולה להיות גדולה מאוד) או כסט של מספר מצומצם של הפעולות האחרונות שהתבצעו. דרך נוחה להצגת הדינמיות של ה-MDP הינה *transition graph* (אוטומט):



Policies and Value functions 25.2.2

מרבית אלגוריתמי ה-*reinforcement learning* כוללים שיערוך של *value functions* המעריכים כמה טוב (=reward עתידי צפוי) ל-*agent* להיות במצב נתון. *reward* עתידי זה תלוי בפעולות שיבחר ה-*agent* לבצע, כלומר ב-*policy*. באופן פורמלי:

הגדרה 25.3 *policy* הינה פונקציית מיפוי בין מצבים לבין ההסתברות לביצוע פעולות מסוימות. אם ה-*agent* פועל לפי פוליסה π בצעד t אזי $\pi(s|t)$ הינה ההסתברות ש- $A_t = a$ אם $S_t = s$

הגדרה 25.4 עבור MDP ופוליסה π נעריך את הפוליסה על ה- MDP באופן הבא:

$$\mathbb{E}_{s_1, s_2, \dots} \left[\sum_{t=0}^{\infty} r(s_t, \pi(s_t)) \right] - Total\ reward \bullet$$

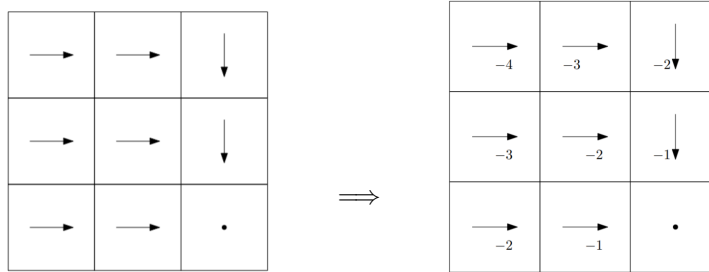
$$\mathbb{E}_{s_1, s_2, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t)) \right] - Discounted\ reward \bullet$$

הגדרה 25.5 ה- $value$ של מצב $s \in \mathcal{S}$ תחת פוליסה π , אשר נסמן $v_{\pi}(s)$, הינו ה- $expected\ return$ כאשר נתחיל ב- s ונפעל על פי π . עבור MDP נתון נגדיר זאת:

$$v_{\pi}(s) = \mathbb{E}_{s_1, s_2, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t)) \mid s_0 = s \right]$$

נקרא ל- v_{π} ה- $state - value\ function$ עבור פוליסה π .

דוגמה-מנחה - חזרה למבוך. נניח שברשותנו פונקציית מעברים דטרמיניסטית ופוליסה. נניח שה- $reward$ הוא -1 לכל המשבצות למעט המשבצת הימנית תחתונה ונניח כי $\gamma = 1$ אזי:



טענה 25.6 לכל פוליסה $\pi : \mathcal{S} \rightarrow \mathcal{A}$ ומצב $s \in \mathcal{S}$ מתקיים כי $v_{\pi}(s) = r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \tau(s, \pi(s))} [v_{\pi}(s')]$

הוכחה: נפתח לפי הגדרה ונקבל:

$$\begin{aligned} v_{\pi}(s) &= \mathbb{E}_{s_1, s_2, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t)) \mid s_0 = s \right] \\ &= r(s, \pi(s)) + \mathbb{E}_{s_1, s_2, \dots} \left[\sum_{t=1}^{\infty} \gamma^t r(s_t, \pi(s_t)) \mid s_0 = s \right] \\ &= r(s, \pi(s)) + \sum_{s' \in \mathcal{S}} p(s' \mid s, \pi(s)) \mathbb{E}_{s_2, \dots} \left[\sum_{t=1}^{\infty} \gamma^t r(s_t, \pi(s_t)) \mid s_0 = s \right] \\ &= r(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} p(s' \mid s, \pi(s)) v_{\pi}(s') \end{aligned}$$

■

באופן דומה נוכל להגדיר זאת עבור בחירת פעולה מסוימת:

הגדרה 25.7 ה- $value$ של בחירת פעולה $a \in \mathcal{A}$ במצב $s \in \mathcal{S}$ תחת הפוליסה π , אשר נסמן $q_\pi(s, a)$, הינו ה- $expected return$ כאשר נתחיל ב- s , נבצע a ונפעל לפי π . עבור MDP נתון נגדיר זאת:

$$q_\pi(s, a) = \mathbb{E}_{s_1, s_2, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \wedge a_0 = a \wedge \forall t > 0 a_t = \pi(s_t) \right]$$

נקרא ל- q_π ה- $action - value function$ עבור פוליסה π .

היחס בין v_π ו- q_π

- v_π כפונקציה של q_π : $\forall s \in \mathcal{S} \quad v_\pi(s) = q_\pi(s, \pi(s))$
- q_π כפונקציה של v_π : $\forall s \in \mathcal{S}, a \in \mathcal{A} \quad q_\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \tau(s, \pi(s))} [v_\pi(s')]$

25.2.3 פוליסות ופונקציות ערך אופטימליות

באופן כללי פתרון בעיית $reinforcement learning$ הינה מציאת $policy$ המניבה $reward$ גבוה בטווח הארוך. עבור MDP נוכל להגדיר זאת באופן הדוק.

הגדרה 25.8 עבור פוליסות $\pi, \pi' : \mathcal{S} \rightarrow \mathcal{A}$ נאמר כי π טובה יותר מ- π' , אם ורק אם $\forall s \in \mathcal{S} \quad v_\pi(s) \geq v_{\pi'}(s)$.

הגדרה 25.9 נגדיר את הפונקציות ערך, פעולה-ערך ופוליסה האופטימליות באופן הבא:

$$\begin{aligned} \text{Opt. value function :} \quad & v_*(s) = \max_{\pi} v_\pi(s) \\ \text{Opt. action - value function :} \quad & q_*(s, a) = \max_{\pi} q_\pi(s, a) \\ \text{Optimal policy :} \quad & \pi_*(s) = \operatorname{argmax}_{\pi} v_\pi(s) \end{aligned}$$

משפט 25.10 לכל MDP מתקיים:

- קיימת π_* אשר טובה או שווה לכל הפוליסות האחרות: $\forall \pi \quad \pi_* \geq \pi$.
- כל הפוליסות האופטימליות משיגות את ה- $optimal value function$: $\forall \pi_* \quad v_{\pi_*}(s) = v_*(s)$.
- על הפוליסות האופטימליות משיגות את ה- $optimal action - value function$: $\forall \pi_* \quad q_{\pi_*}(s, a) = q_*(s, a)$.

היחס בין v_* ו- q_*

• v_* כפונקציה של q_* : $\forall s \in \mathcal{S} \quad v_*(s) = \max_{a \in \mathcal{A}} q_*(s, a)$.

$$\begin{aligned} q_*(s, a) &= q_{\pi_*}(s, a) \\ &= r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi(s))} [v_{\pi_*}(s')] \quad \bullet \text{ } q_{\pi_*} \text{ כפונקציה של } v_{\pi_*} \\ &= r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi(s))} [v_*(s')] \end{aligned}$$

• π_* כפונקציה של q_* : $\pi_*(s) = \operatorname{argmax}_{a \in \mathcal{A}} q_*(s, a)$.

מסקנה 25.11 הפוליסה האופטימלית π_* הינה דטרמיניסטית.

משפט 25.12 (Bellman's Equation) הפוליסה π הינה אופטימלית אם ורק אם מקיימת את משוואת האופטימליות של

$$v_*(s) = \max_{a \in \mathcal{A}} \left[r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi(s))} [v_*(s')] \right] \quad \text{:Bellman}$$

הוכחה:

$$\begin{aligned} v_*(s) &= \max_{a \in \mathcal{A}} q_*(s, a) \\ &= \max_{a \in \mathcal{A}} q_{\pi_*}(s, a) \\ &= \max_{a \in \mathcal{A}} \left[r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi(s))} [v_*(s')] \right] \end{aligned}$$

■

25.3 אלגוריתמים ל-MDP

כאשר אנו באים לפתור בעיות *reinforcement learning* אנו מניחים כי ה-MDP סופי. כלומר $\mathcal{S}, \mathcal{A}, \mathcal{R}$ קבוצות סופיות וכי הדינמיקה שלהם ניתנת לנו על ידי ההסתברויות של $p(s', r | s, a)$ כפי שכבר ראינו. במקרים רבים פתרון בעיית הלמידה נעשה בעזרת תכנון דינמי (*Dynamic Programming*) על מנת למדל את המבנה והחיפוש של פוליסות טובות. ואכן, באופן כללי כאשר נרצה לפתור MDP נשתמש ב-*value functions* כדי למדל את מבנה החיפוש ולמציאת הפוליסה האופטימלית.

בחלק זה נראה כיצד לחשב את פונקציית ה- $state - value$: v_π . בהינתן קלט פונקציית הסתברות מעברים $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ (עבור מצב, הפעולה שנבצע במצב והמצב אליו נגיע בעקבות כך), פונקציה לצפי $reward$ $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ ופונקציית פוליסה $\pi : \mathcal{S} \rightarrow \mathcal{A}$ נרצה להחזיר $v_\pi : \mathcal{S} \rightarrow \mathbb{R}$.

Initialize $\forall s \in \mathcal{S} \ v_0(s) = 0$

For $t = 1, 2, \dots$

$$\forall s \in \mathcal{S} \ v_{t+1}(s) = \sum_{s' \in \mathcal{S}} p(s' | s, \pi(s)) \cdot [r(s, \pi(s)) + \gamma v_t(s')]$$

נבין כעת את החישוב שמבצעים. עבור $s \in \mathcal{S}$ נבדוק לכל $s' \in \mathcal{S}$ מהו הערך שמצב זה מוסיף לפתרון ונסכום באופן ממושקל, לפי ההסתברות של s' להתקיים (על פי המצב s והפעולה שנבצע - $p(s' | s, \pi(s))$, את הערך הצפוי אילו נבחר בכיוון זה וערכו של מצב זה. אלגוריתם שכזה, אשר בכל איטרציה מעדכן את הערך הצפוי של $v_t(s)$ בסופו של דבר מתכנס לערך האמיתי של מצב זה, אולם היות ולא נוכל לתת לאלגוריתם לרוץ עבור $t \rightarrow \infty$ נרצה להחליט באיזה שלב נעצור. לשם כך הטענה הבאה (בכל איטרציה של האלגוריתם מצליחים לשפר את הערך אליו שואפים באופן אקספוננציאלי) וכן הפסאודו-קוד המלא אחרי הטענה.

טענה 25.13 עבור הרצת האלגוריתם *Policy Evaluation* במשך t איטרציות מתקיים $\|v_t - v_\pi\|_\infty \leq \gamma^t \|v_0 - v_\pi\|_\infty$.

הוכחה: לכל $s \in \mathcal{S}$ מתקיים:

$$\begin{aligned} |v_{t+1}(s) - v_\pi(s)| &= \left| \sum_{s' \in \mathcal{S}} p(s' | s, \pi(s)) \cdot [r(s, \pi(s)) + \gamma v_t(s')] \right| - \left| \sum_{s' \in \mathcal{S}} p(s' | s, \pi(s)) \cdot [r(s, \pi(s)) + \gamma v_\pi(s')] \right| \\ &= \gamma \left| \sum_{s' \in \mathcal{S}} p(s' | s, \pi(s)) \cdot [v_t(s') - v_\pi(s')] \right| \\ &\leq \gamma \sum_{s' \in \mathcal{S}} p(s' | s, \pi(s)) \cdot |v_t(s') - v_\pi(s')| \end{aligned}$$

■

Iterative policy evaluation

```
Input  $\pi$ , the policy to be evaluated
Initialize an array  $V(s) = 0$ , for all  $s \in \mathcal{S}^+$ 
Repeat
   $\Delta \leftarrow 0$ 
  For each  $s \in \mathcal{S}$ :
     $v \leftarrow V(s)$ 
     $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
until  $\Delta < \theta$  (a small positive number)
Output  $V \approx v_\pi$ 
```

(מכאן מעט הרחבה על סמך הספר) למעשה המטרה שבחישוב ה-*value function* של הפוליסי היא כדי להצליח למצוא פוליסות טובות יותר. כעת אחרי שחישבנו את הערך הצפוי (כולל התוספת האיטרטיבית הצפויה) של בחירת $a \in \mathcal{A}$ עבור $s \in \mathcal{S}$ כלשהו נוכל לראות האם להחליף לפעולה אחרת ולקבל בסופו של דבר ערך גבוה יותר. לתהליך זה קוראים *policy improvement* (אשר לא נכנסנו אליו בקורס) ונוכל להגדיר פוליסה אשר תפעל באופן חמדן ותבחר להתקדם לפי הערך האופטימלי של $q_\pi(s, a)$:

$$\begin{aligned}\pi'(s) &= \underset{a}{\operatorname{argmax}} q_\pi(s, a) \\ &= \underset{a}{\operatorname{argmax}} \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \underset{a}{\operatorname{argmax}} \sum_{s',r} p(s',r \mid s, a) \cdot [r + \gamma v_\pi(s')]\end{aligned}$$

לאחר השלב של *policy evaluation* ניתן לעשות את ה-*policy improvement* ולחזור על שתי הפעולות הללו עד אשר מגיעים ל-*policy* אופטימלי.

Value Iteration 25.3.2

במקרה זה נרצה לקבל את המצבים והפעולות ולהיות מסוגלים למצוא מה ה-*policy* האופטימלי. אלגוריתם ה-*value iteration* פועל דומה ל-*policy iteration* אך בכל איטרציה בוחר את ה-*action* הממקסם את הרווח העתידי. נפעיל את האלגוריתם שוב ושוב ונתכנס לערך האופטימלי. בדומה ל-*policy iteration* גם כאן קיימת טענה שמראה התקרבות אקספוננציאלית לערך האמיתי כתלות ב- t .

$$\begin{aligned}&\text{Initialize } \forall s \in \mathcal{S} \ v_0(s) = 0 \\&\text{For } t = 1, 2, \dots \\&\forall s \in \mathcal{S} \ v_{t+1}(s) = \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s' \mid s, \pi(s)) \cdot [r(s, \pi(s)) + \gamma v_t(s')]\end{aligned}$$

```

Value iteration
Initialize array  $V$  arbitrarily (e.g.,  $V(s) = 0$  for all  $s \in \mathcal{S}$ )

Repeat
   $\Delta \leftarrow 0$ 
  For each  $s \in \mathcal{S}$ :
     $v \leftarrow V(s)$ 
     $V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
until  $\Delta < \theta$  (a small positive number)

Output a deterministic policy,  $\pi \approx \pi_*$ , such that
 $\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$ 

```

הערה 25.14 נשים לב שב- $\Delta < \theta$ אנו בעצם מגדירים שהאלגוריתם ממשיך איטרטיבית כל עוד Δ לא התכנס (כלומר עד שהשיפור יהיה מספיק קטן, תחת ההנחה שמכאן התקרבונו מספיק לערך האמיתי).

הערה 25.15 להסבר יותר מפורט על ההבדל בין *value iteration* ו-*policy iteration* ניתן לקרוא את השאלת *stackoverflow*

הבאה: What is the difference between value iteration and policy iteration?

(או בקישור <https://stackoverflow.com/q/37370015/6400526>)

Learning 25.3.3

עד כה עסקנו בבעיית התכנון אשר בה ידענו מהם ה-*rewards* ופונקציית ה-*transition probability* ויכולנו להסיק מה הדרך המועדפת. כאשר איננו יודעים פונקציות אלה הקלט לאלגוריתם הלמידה הוא התנסויות. במקרה זה הם זוגות של *state, action* נוכחיים ושל תוצאה מביצוע הנוכחיים (הגדרה נוספת היא סדרה של זוגות עוקבים). ישנם שני סוגי למידה שכאלה:

- *On policy* - הפעולות מתבצעות על ידי ה-*agent* עצמו והוא שמחליט כיצד לפעול.
- *Off policy* - הגישה בה "מתבוננים מהצד". נותנים לפעולות להתבצע אך לא משפיעים עליהן, מודדים תוצאות ומסיקים מסקנות.

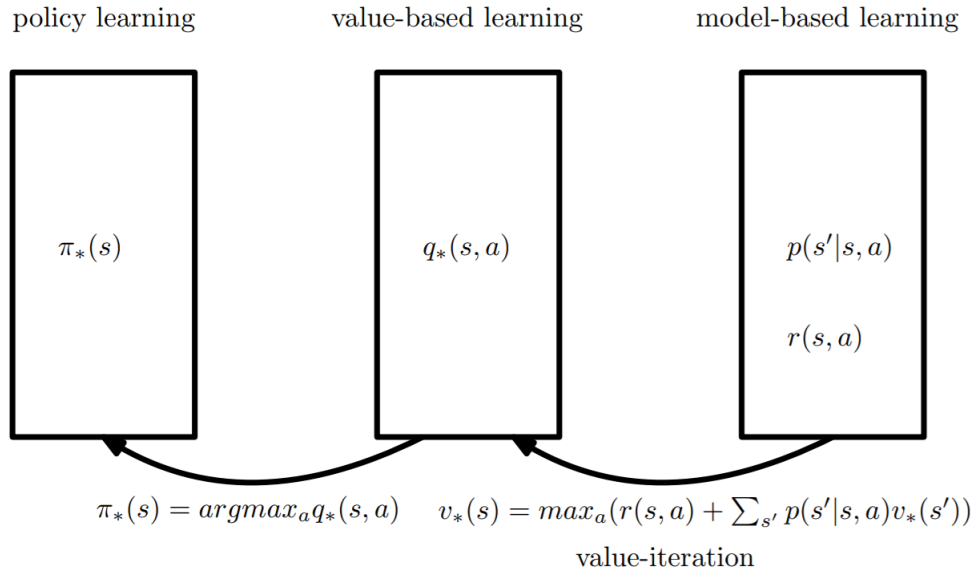
להבדלים בין השניים ניתן לקרוא כאן What is the difference between off-policy and on-policy learning? (או בקישור <https://stats.stackexchange.com/q/184657>)

כאשר נשאלת השאלה מה אנו רוצים ללמוד יש כמה גישות:

1. לימוד מבוסס מודל - נלמד תחילה את ה-*MDP*. נלמד את הרווחים ואת ה-*transition probabilities*. בהתאם לזאת נוכל ללמוד איך להתנהג אחר כך. מכאן כדי להסיק את ה-*value function* (בהינתן שכעת אנו יודעים את המודל) נוכל להשתמש ב-*value iteration* שראינו.
2. לימוד מבוסס value - ישנם אלגוריתמים ללמידת ה-*q-function*. בהינתן ה-*q-function* נמצא את הפוליסה האופטימלית על ידי כך שניקח את ה-*actions* האופטימליים.

3. לימוד מבוסס פוליסה - אלגוריתמים ללמידה ישירה של ה-policy האופטימלי.

אלגוריתמים של שני הסעיפים האחרונים נקראים *model free algorithms*.



אלגוריתם Q – Learning

באלגוריתם הבא אנו מקבלים כקלט התנסויות ונעדכן את ה-*q function* בהתאם. צורת העדכון, אשר דומה לצורה ב-GD, לוקחת את הערך הנוכחי עם תיקון קטן לפי ההפרש בין הערך החדש לקודם. אז זה נכפיל בפקטור η אשר משמש כ-*learning rate*. נשים לב כי את הערך החדש איננו יודעים (שכן איננו יודעים את פונקציית ההסתברות) ולכן נשערכו בהתאם לכן שניקח את הפעולה שמביאה למקסימום.

Initialize $\forall s \in \mathcal{S}, a \in \mathcal{A} \ q_0(s, a) = 0$

For each experience (s_t, a_t, r_t, s_{t+1})

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \eta_t \left(q_t(s_t, a_t) - \left(r_t + \gamma \max_a q_t(s_{t+1}, a) \right) \right)$$

משפט 25.16 אם כל (s, a) מתקיים *infinitely often* ו- $\eta_t = \frac{1}{t}$ אזי $\lim_{t \rightarrow \infty} q_t(s, a) = q_*(s, a)$.

נשים לב כי אלגוריתם זה הינו *off policy* שכן קיבלנו התנסויות כלשהן ועדכנו מסקנות בהתאם.

אם נרצה לבחור אלגוריתם שהינו *on policy* אשר בו נשפיע על איזה *actions* יתבצעו, נתקל בדילמה שהזכרנו בתחילת הפרק של *exploration vs. exploitation* - האם נרצה שהאלגוריתם יקח את ה-policy הכי טוב (ה-action שיביא

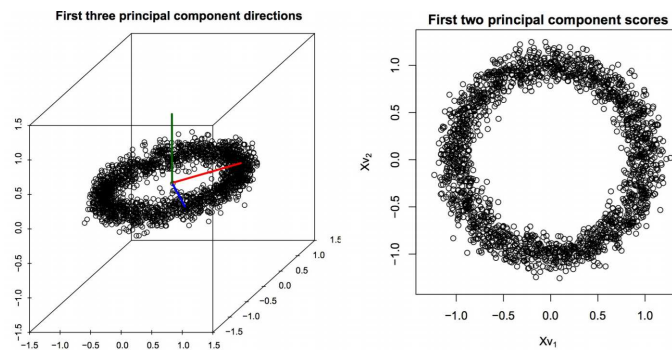
למקסימום) או לא ואולי נמצא כך מסלולים חדשים טובים יותר שלא ידענו עליהם. שיטה גנרית לעשות זאת ($\epsilon_t - greedy$) היא שעבור כל מצב בהסתברות $1 - \epsilon_t$ נבצע $exploitation$, כלומר נבחר את ה- $action$ הממקסם, ובהסתברות ϵ_t נבצע $exploration$, כלומר נבחר $action$ אקראי.

נשים לב שכאשר מספר המצבים גודל ההסתברות לעשות $exploration$ היא מאוד נמוכה. על כן נשלב בין $reinforcement learning$ ורשתות נוירונים.

26 Dimensionality Reduction

בפרק זה נעסוק בהתמודדות עם $feature space$ ממימד גבוה מדי. מעבר לקושי החישובי שבמספר $features$ גדול קשה לנו לחקור את ה- $data$ ולנתח את התוצאות. בשלב זה ראינו מספר דרכים כמו $lasso$ ו- $subset selection$ להוריד $features$.

במקרים רבים, על אף הריבוי ב- $features$, יש יתירות ($redundancy$) גבוהה - כלומר אין לנו צורך בכל ה- $features$ כדי להצליח לתייג את המידע. על כן יתכן ונוכל לקפל למימד נמוך יותר. תהליך ה- $dimensionality reduction$ לוקח $data$ ממימד גבוה (המרחב האמביינטי) וממפה אותו למרחב חדש, שמימדו קטן יותר (המרחב האינטרניזי), אשר עדין מגלם את אותו עולם ראשוני. במקרים רבים פעולת הורדת המימד הינה שלב מקדים ללמידה.



בתמונה זו אנו רואים $data$ שנמצא במרחב התלת מימדי וטרנספורמציה למרחב הדו מימדי שעדין מייצגת נכונה את ה- $data$. כלומר יכולנו להוריד את המימד הדגימות וזאת מבלי לאבד יותר מדי מהוריאביליות של ה- $data$.

26.1 Linear Dimension Reduction and PCA

ברוב המקרים, עבור \mathcal{X} הנמצא ב- \mathbb{R}^d מתקיים שבפועל $\mathcal{X} \subseteq \mathbb{R}^r \subset \mathbb{R}^d$. כלומר \mathcal{X} הינו תת מרחב ממימד קטן יותר r . נרצה אפוא לחפש העתקה לינארית $w : \mathbb{R}^d \rightarrow \mathbb{R}^r$ כך שנוכל לייצג את המידע שלנו במימד נמוך יותר. העתקה זו היא למעשה $x \mapsto Wx$ כאשר $x \in \mathbb{R}^d$ ו- $Wx \in \mathbb{R}^r$. יהיו $x_1, \dots, x_m \in \mathbb{R}^d$ דגימות אזי נחפש מטריצות $W \in \mathbb{R}^{r \times d}$ ו- $U \in \mathbb{R}^{d \times r}$ כך ש:

• $x \mapsto Wx$ הינה הטרנספורמציה כך ש- Wx הינו הייצוג של x במימד הנמוך יותר r .

• נסמן $y = Wx$ אזי נוכל לחשב את $\tilde{x} = Uy$ אשר הוא שיחזור מקורב של x במימד הגבוה d .

הגדרה 26.1 ב-*Principal Component Analysis (PCA)* נחפש מטריצת כיווץ W ומטריצת שחזור U כל שסכום ריבועי המרחקים בין הוקטור המקורי והמשוחזר מינימלי:

$$\operatorname{argmin}_{W \in \mathbb{R}^{r \times d}, U \in \mathbb{R}^{d \times r}} \sum_{i=1}^m \|x_i - UWx_i\|_2^2$$

הערה 26.2 מטרת המטריצה U למעשה היא לספק כלי השוואה בין המרחב אליו הגענו לאחר הורדת המימד לבין המרחב המקורי.

הערה 26.3 נשים לב כי איננו יודעים את r האמיתי של ה- $data$ ועל כן בפועל אנו מחפשים למעשה k כלשהו שבתקווה הינו r .

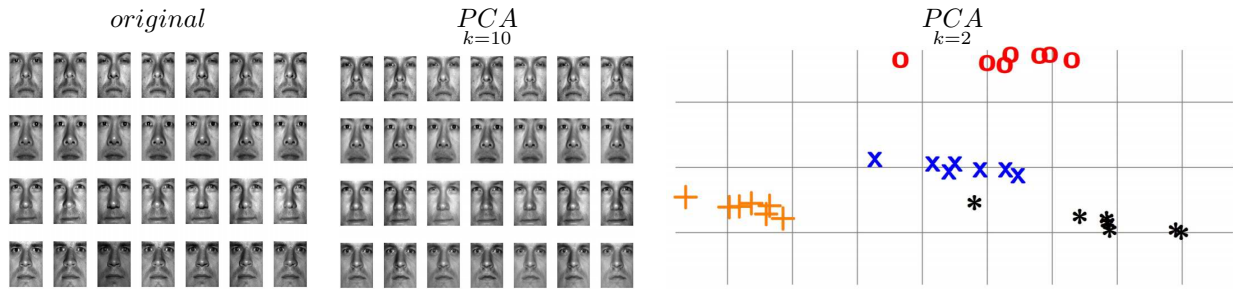
משפט 26.4 עבור הורדת מימד עם PCA ממימד d למימד k כלשהו ומטריצות $U \in \mathbb{R}^{d \times k}$, $W \in \mathbb{R}^{k \times d}$ כאשר עמודות U הן k הוקטורים העצמיים הגדולים של מטריצת ה- $covariance$ S אזי $W = U^T$

$$S = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^T$$

נבחין כי במטריצה S (בנוסף למרכז - ההחסרה של הערך הממוצע בכל נקודה) מתקיים שהיא מטריצה סימטרית עם ערכים עצמיים אי שליליים שעמודותיה בסיס אורתונורמלי לתת המרחב. למעשה עבור $\lambda_1, \dots, \lambda_d \geq 0$ אם נניח בה"כ $\lambda_1 \geq \dots \geq \lambda_d \geq 0$, מתקיים כי הוקטורים העצמיים השייכים לערכים העצמיים בעלי ערכים קטנים מייצגים מעט מאוד וריאביליות במידע. לעומת זאת הוקטורים העצמיים השייכים לעצמיים גדולים מייצגים וריאביליות גבוהה של המידע בכיוון של הוקטור העצמי. על כן אם ניקח את k הוקטורים העצמיים עם הערכים העצמיים הכי גדולים נקבל את תת המרחב ממימד k שמייצג הכי טוב את הוריאביליות של המידע המקורי - ובכל זאת במימד נמוך יותר.

למעשה עבור u_1, \dots, u_d הוקטורים העצמיים השייכים לערכים העצמיים $\lambda_1, \dots, \lambda_d$ ועבור המטריצה X עם m דגימות ו- d *features*, כדי לקבל הורדת מימד למימד k ניקח את $Xu_1u_1^T, \dots, Xu_ku_k^T$ עבור k כרצוננו.

דוגמה-מנחה: עבור הבעיה של זיהוי פנים, בהינתן ה- $dataset$ של התמונות השונות, היות ואין באמת שוני רב בין התמונות, למעט האדם השונה, נוכל להפעיל PCA עם ערך k כלשהו (בדוגמה $k = 10$) ולמעשה לקבל $data$ בבמימד נמוך יותר שעדין מתאר טוב מספיק את ה- $data$ המקורי. נמשיך ונראה כי עבור $k = 2$, ההטלה של התמונות מממד 50×50 למימד 2 מספקת לנו אזורים שונים שבכל אזור למעשה קיבלנו את אוסף הנקודות, התמונות, שמייצגות אדם נפרד.



Clusters – k – means 27

במשך הקורס עסקנו ב-*supervised learning* שבו לדגימות שברשותנו יש תיוגים. במקרה של *unsupervised learning* המידע איננו מכיל תיוגים. שתי דוגמאות לכך הן *clustering* ו-*anomaly detection*. באופן אינטואיטיבי קליסטור (*clustering*) הינה הפעולה של קיבוץ המידע לסטים כך שפריטים בעלי אופי דומה יהיו באותה קבוצה, ופריטים בעלי אופי שונה יהיו בקבוצות שונות. לדוגמה תחת ביולוגיה חישובית ננסה לקלסטור קבוצות של גנים בהתבסס על דמיון ברמות הביטוי שלהם תחת ניסויים שונים. באופן יותר פורמלי עבור:

קלט: יהי \mathcal{X} סט המידע ו- $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ פונקציית מרחק מעל \mathcal{X} המקיימת סימטריות ו- $d(x, x) = 0 \forall x \in \mathcal{X}$. לעיתים נדרוש שתקיים גם את אי שיויון המשולש. לעיתים נגדיר את פונקציית המרחק להחזיר ערכים $[0, 1]$ ו- $d(x, x) = 1$. ישנם אלגוריתמים שגם יבקשו $k \in \mathbb{N}$ הקובע את מספר ה-*clusters*.

פלט: חלוקה של הדומיין \mathcal{X} . כלומר, $\bigcup_{i=1}^k C_i = \mathcal{X}$ s.t. $C = (C_1, \dots, C_k)$. נזכור כי בחלוקה הכוונה היא $\forall i \neq j, C_i \cap C_j = \emptyset$. לעיתים ב-*soft clustering* נחזיר וקטור הסתברויות לשייכות ל-*cluster*. כלומר לכל $x \in \mathcal{X}$ נחזיר $(p_1(x), \dots, p_k(x))$ s.t. $\forall i \in [k] p_i(x) = \mathbb{P}(x \in C_i)$

האלגוריתמים השונים מגדירים פונקציה לחישוב מרחק בין *clusters* (לדוגמה זוג הנקודות הכי קרובות בין השני *clusters*, הזוג הכי רחוק או ממוצע המרחקים) וכן פונקציית *cost* עבור *cluster*. על כן נחפש חלוקה, *clusters*, שיביאו את העלות למינימום. כלומר $G : (\mathcal{X}, d) \times C \rightarrow \mathbb{R}_+$ עבור $C = (C_1, \dots, C_k)$ חלוקה של הדומיין (הסימון G נובע מ-*goal*). בעיות מינימיזציה אלה הן *NP-hard* וחלקן אפילו קשות לקירוב ועל כן במקרים רבים נבקש גם את הפרמטר k עבור כמות ה-*clusters* הרצויה (במקום שהאלגוריתם ימצא לבד).

כאשר דנים באלגוריתם *k-means*, מחפשים חלוקה C_1, \dots, C_k כאשר C_i מיוצג על ידי הצנטרואיד μ_i (*centroid*). מניחים כי \mathcal{X} מוכל במטריקה רחבה יותר \mathcal{X}' אשר נקודות בצנטרואיד נמצאות ב- \mathcal{X}' . כלומר נקודות הצנטרואידים אינם בהכרח נמצאים ב- \mathcal{X} . על כן ב-*k-means* מודדים את המרחק הריבועי של כל דגימה מהצנטרואיד של ה-*cluster*:

$$\mu_i(C_i) \stackrel{\text{def.}}{=} \underset{\mu \in \mathcal{X}'}{\operatorname{argmin}} \sum_{x \in C_i} d(x, \mu)^2$$

ומכאן כי עבור דומיין \mathcal{X} ופונקציית מרחק בין דגימות d :

$$G_{k\text{-means}}((\mathcal{X}, d), (C_1, \dots, C_k)) = \sum_{i=1}^k \sum_{x \in C_i} d(x, \mu_i(C_i))^2$$

\Updownarrow

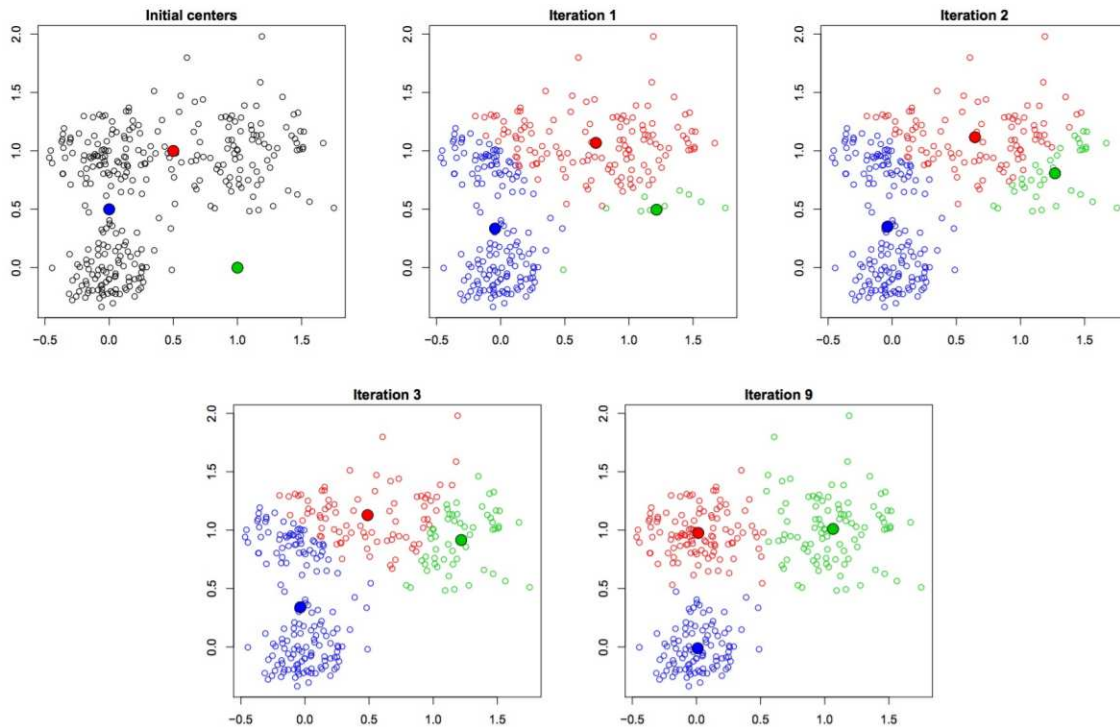
$$G_{k\text{-means}}((\mathcal{X}, d), (C_1, \dots, C_k)) = \min_{\mu_1, \dots, \mu_k \in \mathcal{X}'} \sum_{i=1}^k \sum_{x \in C_i} d(x, \mu_i)^2$$

כפי שכבר תיארנו בעיות אלה הן $NP - hard$ ולכן אחת היוריסטיקות שבשימוש היא *Lloyd's algorithm* אשר משתמשת ב־*alternating minimization*:

1. עבור סט דגימות x_1, \dots, x_m ומספר $clusters$ רצוי $k \in \mathbb{N}$.
2. נבחר באופן מקרי μ_1, \dots, μ_k צנטרואידים.
3. נפעל באופן איטרטיבי עד להתכנסות (עד אשר ההבדל בצנטרואידים באיטרציה t לאיטרציה $t+1$ קטנה מסף שנחליט).

(א) נגדיר C_i קבוצת הנקודות שהכי קרובות לצנטרואיד μ_i מכל צנטרואיד אחר.

(ב) נעדכן את הצנטרואיד μ_i להיות מרכז הכובד של הקבוצה: $\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$



הערה 27.1 קיימים משפטים המוכיחים כי בכל איטרציה השונות בתוך *cluster* יורדת וכי לבסוף יש התכנסות.

הערה 27.2 חשוב לשים לב שהערכים הראשונים שהצטרואידים מקבלים משפיעים דרמטית על ה-*clusters* שנקבל. כלומר ה-*variance* של אלגוריתם זה גדול. נרצה למשל להריץ מספר פעמים ולראות מה האופציה שיצאה טובה מבחינתנו. בנוסף, נזכור גם שאנחנו קובעים את ה-*k* עבור מספר ה-*clusters*. לא מובטח שה-*data* אכן ניתן להפרדה טובה (והגיונית על פי מה שה-*data* מייצג) ל-*k* קבוצות.

