

דה-קורלציה: מבוצעת כיוון שניתן להוריד את השונות רק עד מידה מסוימת. למשל בעצים נבחר פרמטר k ובכל ענף נבחר k פרמטרים שרק מתוכם ניתן לפצל.

Bias & Variance - Bagging		
Bias	Doesn't change	Bias
Low Var	Big(∞)	Small
High Var	Small	Big(∞)

Random Forest שימוש ב-*Bagging* כדי לשפר את הביצועים של עצי החלטה.

Random Forest Algorithm:

- מאתחלים T עצים בעומק R באמצעות m_{min} דגימות לכל עץ $t \in \{1, \dots, T\}$:
- מייצרים דגימת $bootstrap$ - S^* מתוך S
 - מאמינים על החלטה עם הדגימה ומקבלים הפיטורה h_{S^*} .
 - כל עוד לא הגענו לעומק מקסימלי או מספר מינימלי של נקודות בקוביה:
 - בוחרים בנצירה אחידה k פיצורים מתוך d הפיצורים
 - בוחרים את הפיצור הטוב מבניהם ומבצעים הפירה לפיו
- נחזיר את $h_S(x) = \text{sign}(\sum_{t=1}^T w_t h_t(x))$

Boosting (הורדת ה-Bias, Variance עולה לאט) ע"י עדכון מרחב המדגם S^t (הוספת משקולות) כך שדגימות בהן טעינו יקבלו משקל גדול יותר (יצור התפלגות D^t). [או שניתן לדגום מחדש מתוך התפלגות D^t , תמיד זמין איך יכול ליצור כפילויות הדגימה]

Adaboost Algorithm:

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

Loop: for t in $[T]$: $h_t = \mathcal{A}(D^t, S)$

Update:

- $\sum_{j=1}^m D_j^{t+1} \cdot 1_{[h_t(x_j)=y_j]} = \frac{1}{2} \Rightarrow w_t = \frac{1}{2} \log(\frac{1}{\epsilon_t} - 1)$ (סיכון אמפירי ממושקל, $\epsilon_t = \sum_{j=1}^m D_j^t \cdot 1_{[h_t(x_j) \neq y_j]}$)
- $D_j^{t+1} = \frac{D_j^t e^{-w_t y_j 1_{[h_t(x_j) \neq y_j]}}}{\sum_{j=1}^m D_j^t e^{-w_t y_j 1_{[h_t(x_j) \neq y_j]}}}$

לשים לב: כאשר נגדיל את T Adaboost- T (כלומר את מספר ההלומים החלישים) אנו נשפר את $\text{error} - \text{Approximation}$ (הבלתאם את $\text{error} - \text{Estimated}$).

חיזוי דגימה חדשה:

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

Bias & Variance - Adaboost			
Low Bias	Big(∞)	Small	High Bias
Low Var	Small	Big(∞)	Var goes up a bit

לומדים חלשים γ -weak-learner: אלוגוריתם למידה \mathcal{A} הוא לומד- γ חלש עבור מחלקת היפותזות \mathcal{H} אם קיימת פונקציה $\eta: (0,1) \rightarrow \mathbb{N}$ כזו: $m_{\mathcal{H}}: (0,1) \rightarrow \mathbb{N}$ שכל $\delta \in (0,1)$, לכל התפלגות \mathcal{D} מעל מרחב מדגם \mathcal{X} , ולכל פונקציית תיוג $f: \mathcal{X} \rightarrow \{\pm 1\}$, אם הנחת הריאליזביליות מתקיימת (ביחס $(\mathcal{H}, \mathcal{D}, f, \cdot)$), אז כאשר נריץ את \mathcal{A} על \mathcal{D} מספר $m \geq m_{\mathcal{H}}(\delta)$ פעמים יחזיר \mathcal{A} האלגוריתם חיזוי הפיטורה $h_S = \mathcal{A}(S)$ כך ש:

$$\Pr_{S \sim \mathcal{D}^m} \left[L_{\mathcal{D},f}(h_S) \leq \frac{1}{2} - \gamma \right] \geq 1 - \delta$$

γ -weak-learnable: מחלקה \mathcal{H} היא למידה- γ חלשה אם קיים \mathcal{A} כ"ל.

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

משפט: T הן קבוצת דגימות. ניח שבכל איטרציה של Adaboost תלומד הבסיסי מחזיר לכל חיזוי (היפוזתה h_t) עבודה הסיכון האמפירי הממושקל מקיים

$$\epsilon_t = \sum_{j=1}^m D_j^t \cdot 1_{[h_t(x_j) \neq y_j]} \leq \frac{1}{2} - \gamma$$

אז כלל הסיכון המפירי (הלא ממושקל) של כלל ההחלטה של Adaboost מקיים:

$$L_S(h_{\text{boost}}) \equiv \frac{1}{m} \sum_{j=1}^m 1_{[h_{\text{boost}}(x_j) \neq y_j]} \leq e^{-2\gamma^2 T}$$

Boosting	Bagging	מקביליות
טורי	במקביל	
מבנה הנתונים	Bootstrap training samples	מבנה הנתונים
De-correlation	מומלץ	
Overfitting	לא קורה	
איזה מודל להשתמש בו	עצים עמוקים	
הטענה	מפחית Variance	מפחית Bias

גולרדציה $\mathcal{A}_S: S \rightarrow h_S$ על \mathcal{X} נגדיר לומד h_S על \mathcal{X} כך ש:

$$h_S = \arg \min_{h \in \mathcal{H}} \mathcal{R}(h) + \lambda \cdot \|h\|_0$$

Regularization Term: \mathcal{R} הנוטה להעניף את הנתונים הרגולציה. עבור $\lambda = 0$ אין התייחסות כלל למורכבות ההיפוזתה שנוצרה. ה-Bias נמוך אך Variance גבוה.

- עבור $\lambda \rightarrow \infty$ נעדיף היפוזתה כמה שיותר פשוטה ללא קשר לפונקציית המטרה

Bias & Variance - Regularization		
Low Bias	0	$\xrightarrow{\lambda} \infty$
Low Var	∞	$\xrightarrow{\lambda} 0$

הרגולרציה של רגריסה לינארית

Ridge Regression

$$\arg \min_{w_0 \in \mathbb{R}, w \in \mathbb{R}^d} \|w_0 \cdot \vec{1} + Xw - y\|_2^2 + \lambda \cdot \|w\|_2^2$$

יש לפתור את המערכת $X^T y = (X^T X + \lambda I)w$ מתקיים $y^T y = U^T V^T y$ כאשר $\lambda > 0$

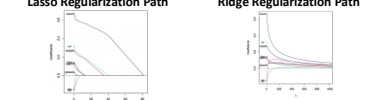
$$S^\lambda = \text{diag}\left(\frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right)$$

Lasso - בעל תכונת דילוח.

$$\arg \min_{w_0 \in \mathbb{R}, w \in \mathbb{R}^d} \|w_0 \cdot \vec{1} + Xw - y\|_2^2 + \lambda \cdot \|w\|_1$$

Best Subset

$$\arg \min_{w_0 \in \mathbb{R}, w \in \mathbb{R}^d} \|w_0 \cdot \vec{1} + Xw - y\|_2^2 + \lambda \cdot \|w\|_0$$



טענה: עבור $w \in \mathbb{R}^d$ ומטריצה $X^T X = I$ (תכנון אורתוגונלי):

- $\hat{w}_\lambda^{\text{ridge}} = \frac{1}{1+\lambda} \hat{w}^{\text{LS}}$, where $\hat{w}^{\text{LS}} = (X^T X)^{-1} X^T y = X^T y$
- $\hat{w}_\lambda^{\text{lasso}} = \eta_\lambda^{\text{soft}}(\hat{w}^{\text{LS}})$ s.t. $\eta_\lambda^{\text{soft}}(x) = \begin{cases} x - \lambda & x \geq \lambda \\ 0 & -\lambda < x < \lambda \\ x + \lambda & -x \geq \lambda \end{cases}$
- $\hat{w}_\lambda^{\text{subset}} = \eta_\lambda^{\text{hard}}(\hat{w}^{\text{LS}})$ s.t. $\eta_\lambda^{\text{hard}}(x) = x \cdot \mathbb{1}_{|x| \geq \lambda}$

Model Selection

k-fold Cross Validation: עבור $k = 1, \dots, k$ נאמן את המודל i במספר על כל הדאטה מלבד החלק i . נחשב את ה-loss של המודל i על i החלק i . נחזיר את הממוצע וסטיית התקן של התוצאות.

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

Bootstrap: בדגום B קבוצות S^b , ונאמן את האלגוריתם עליהם. בכדי לבחון את הביצועים שנשתמש ב- S^b $T^b = S \setminus S^b$ ונמצע במקדום.

טעויות בבחירת מודל

- Over-estimating generalization error** – כאשר מאמינים $\frac{k-1}{k}$ כמות לא מספקת של דגימות, אז לפי PAC נעריך בחסר את השיגור.
- Under-estimating generalization error** – כאשר אנו מעריכים מודל על פני נתונים המודל אומן (מצב של Overfitting), או כאשר על הדאטה בוצע עיטוי המתאים למודל.

השיגור בעת שימוש בולידציה

$$L_D(h_S) = L_D(h_S) - L_P(h_S) + L_P(h_S) - L_S(h_S) + L_S(h_S) + \frac{L_S(h_S)}{C}$$

A. שיגור הכללה. ניתנת לחסימה תחת ההנחה שפונקציית ה-loss סתומה. B. עבור C גדול $C-1$ ככל הנראה Overfitting . C. עבור C גדול כלל הנראה Underfitting.

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

$$\Pr \left[|L_P(h) - L_D(h)| \leq \sqrt{\frac{\ln(2/\delta)}{m}} \right] \geq 1 - \delta$$

עבור $\mathcal{H}_1 \subseteq \dots \subseteq \mathcal{H}_k$ סופיים.

Standard Method: עבור $h^* \in \text{ERM}_{\mathcal{H}_k}(S_{all})$:

$$\Pr \left[L_D(h^*) \leq \min_{h \in \mathcal{H}_k} L_D(h) + \sqrt{\frac{2 \ln(2k/\delta)}{m}} \right] \geq 1 - \delta$$

Model Selection: עבור $\mathcal{H}_1 \subseteq \dots \subseteq \mathcal{H}_k$ וכל $h^* \in \arg \min_{h \in \mathcal{H}_k} L_D(h)$:

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

Unsupervised Learning

דוגמאות לשימוש: 1. חשיפה של מבנה במידת נמוך. 2. קלאסיפיקציה. 3. זיהוי אנומליות.

PCA

$$\arg \min_{U \in \mathbb{R}^{d \times k}, U^T U = I} \sum_{i=1}^k \|x_i - U U^T x_i\|^2$$

משפט: תהי $X = \sum_{i=1}^m x_i x_i^T$ ויהי u_1, \dots, u_m הו"ע של X (מסודרים בסדר יורד), אז $U^T = [u_1 \dots u_k]$ פתרון לבעיית PCA.

- $P = \sum_{i=1}^m u_i u_i^T$ נקראת מטריצת הטלה, היא סימטרית והיא קרוב טוב ביותר ל- X בתת-מרחב: $\|x - P x\|_2 \geq \|x - P x\|_2 \forall x \in \text{span}\{u_1, \dots, u_k\}$
- היא שונה מ- I בו מטילים את הערך γ לבדו.

ע"י הזהה של הנקודות $\bar{x}_i = x_i - \bar{x}$ אנו מקבלים A את מטריצה השונות: $A = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^T$

הגדרות: תהי מטריצת השונות של דגימות אימון x_1, \dots, x_m כנומדר לעיל, ויהי u_1, \dots, u_d הו"ע של מטריצת השונות המתאימות לע"ע $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$.

המספר λ_1 נקרא **הערך המוביל** ה- i של x_1, \dots, x_m .
הקוטור u_i נקרא **הקוטור המוביל** ה- i של x_1, \dots, x_m .

הורדת מימד: לוקחים את $(x - \bar{x})(x - \bar{x})^T = \sum_{i=1}^m x_i x_i^T$ (המטריצה $Covariance(X)$, מחשבים את הו"ע שלה ולוקחים את k הגדולים מבניהם (קוטורים מובילים). מדברים אותם כעמודות ומטריצה $U^T X$ כעת $U^T X$ פותרת את PCA – הטלה לתת המרחב הנוכח מממד k (אלו נקודות אינטרדיטיות – נקודות בתוך תת המרחב \mathcal{C}_k הפעלנו מממד k). $U^T X$ נותנת את הקואורדינטות של הנקודה לאורך הטלה אבל **במרחב המוקטן** (מממד k).

לאחר מכן כל קוטור במרחב ניתן ליצוג ע"י הקירוב $\bar{x} + \sum_{i=1}^k a_i u_i$.
אם יש יותר דגימות מפיצורים, חישוב A הוא כבד $O(d^3)$ ולכן נרצה להשתמש באלגוריתם הבא:

PCA Algorithm:

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

Eval:

- if $m > d$:
 - $A = X^T X$
 - u_1, \dots, u_k eigenvectors of A
- else ($m \leq d$):
 - $B = X X^T$
 - v_1, \dots, v_k eigenvectors of B
 - denote $u_i = \frac{1}{\|X^T v_i\|} X^T v_i$
- return u_1, \dots, u_k

כשכל לבחור את k ניתן לייצר **Scree Plot** (גרף בו צופים את ה"ערך בסדר יורד).

Clustering

נגדיר חלוקה של המידע x_1, \dots, x_m ל- k מחלקות C_1, \dots, C_k . נגדיר את פונקציית המחיר (ע"פ פונקציית מרחק d) ע"י:

$$G(C_1, \dots, C_k) = \min_{\mu_1, \dots, \mu_k} \sum_{j=1}^k \sum_{x \in C_j} d(x, \mu_j)$$

Centroid: (של C_j) הערך μ_j

$$\arg \min_{\mu_j} \sum_{x \in C_j} d(x, \mu_j)$$

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

k-means Algorithm:

Input: x_1, \dots, x_m and $k \in \mathbb{N}$

Step 0: choose initial μ_1, \dots, μ_k

Until convergence:

- set C_j to be the points x_i closer to μ_j than to any other centroid.
- update μ_j to centroid of C_j : $\mu_j = \frac{1}{|C_j|} \sum_{x \in C_j} x$

החלוקה המושגת μ_1, \dots, μ_k נקראת **Voronoi cells**.

תכונות: השונות בכל C_j יורדת בכל איטרציה. האלגוריתם תמיד מתכנס. האלגוריתם מתכנס למינימום מקומי (בהתאם לנקודות ההתחלה).

Bias & Variance - Clustering		
Low Bias	Big	Small
Low Var	Small	Big

Spectral Clustering

אלגוריתם שמבצע צעדים מקדימים חכמים לפני שהוא מפעיל k -means, נרצה להתעלם ממרחקים גדולים מדי, כלומר נחשוב עליהם בתור אינסוף ונחשב רק במרחקים קטנים. ואז נבנה גרף מיטרי ממושקל שהקודקודים שלו הם מהדגימות שלנו, וזהיה קשת בין שני קודקודים אם הם במרחק ביניהם קטן (כך נבוא ליידי ביטוי להתעלמות ממרחקים גדולים) נבניס ברכיבי הקשירות של הגרף ולפיו נקבע את הקלאסטרים.

נגדיר את מטריצת הסמכויות $A_{i,j} = \exp(-\frac{\|x_i - x_j\|^2}{\epsilon})$ (שמגלמת שכל שהדגימות יותר דומות האפיינטי שלהן יותר גדול, אם הוא 1 הן זהות). נמצא את הו"ע של A $L = D^{-1} A$ (שנקראת גרף לפליסיאן), כאשר $A_{i,j} = \sum_{k=1}^m f(C_k(x_i), C_k(x_j))$. הם ייצגו את הקלאסטרים, וכל דגימה i תיוצג במימד k ע"י הקוארדינטות i -ה שלהם. כעת נפעיל k -means k במימד k .

Kernel methods

סוג של בעיה הפוכה מ-PCA – אנחנו רוצים להטיל למרחב \mathbb{R}^k יותר נוח להפיק מידע נוסף מהדגימות. נפתור בעיות אופטימיזציה המצורה:

$$w^* = \arg \min_{w \in \mathbb{R}^k} \|f((w|\psi(x_1)), \dots, (w|\psi(x_m)))\|^2 + R(\|w\|_2^2)$$

לפי משפט ההייצוג קיים $\alpha \in \mathbb{R}^m$ כך ש- $w^* = \sum_{i=1}^m \alpha_i \psi(x_i)$. תהי $G \in \mathbb{R}^{m \times m}$ המוגדרת ע"י: $G_{i,j} = \langle \psi(x_i), \psi(x_j) \rangle$. הבעיה הנ"ל שקולה ל: $\arg \min_{\alpha \in \mathbb{R}^m} f(G\alpha) + R(\alpha^T G \alpha)$ והוחזיו מתבצע ע"י: $k_i = \langle \psi(x_i), \psi(x_i) \rangle$, כאשר $y(x) = (w^*|\psi(x)) = \alpha^T k$. פונקציה סימטרית k תיקרא **PSD-kernel** אם לכל $m \in \mathbb{N}$ ו- $x_1, \dots, x_m \in \mathcal{X}$ ולכל $k_i = \langle \psi(x_i), \psi(x_i) \rangle$ המטריצה $K(x_1, \dots, x_m) = \sum_{i,j=1}^m k_{ij}$ היא PSD. לכן ע"מ להראות שהעניקה k כלשהי היא פונקציית קרנל נראה 1. מטריצת הגרף המצורה אינה היא PSD. 2. אז קיימת פונקציה ψ שעברה היא ממשתת מכפלה פנימית במרחב החדש.

תנאי מרחב: תהי $\mathcal{X} \rightarrow \mathcal{F}$ פונקציה ψ כאשר \mathcal{F} מרחב הילברט, אזי $\mathcal{X} \rightarrow \mathcal{F} \rightarrow \mathbb{R}$ ממשתת ב- k ממשתת ב- k מכפלה פנימית אם היא PSD-kernel.

יצירת קרנלים חדשים:

יהי $K_1(x, x'), K_2(x, x')$ פונקציות קרנל חוקיות, כל הבאים הם קרנלים חוקיים:

- $f(x, x') = K_1(x, x') f(x')$ לכל פונקציה f .
- $\alpha K_1(x, x') + \beta K_2(x, x')$ עבור $\alpha, \beta \geq 0$.
- $K(x, x') = K_1(x, x') K_2(x, x')$.

דוגמאות לחדשים:

בעיית ה ridge regression מקיימת: $y(x) = k^T(G + \lambda I)^{-1} y$

פונקציית קרנל בעבור ההאמתה פולינומיליאלית מדרגה לכל היותר k מתקבלת ע"י: $k(x, x') = (1 + \langle x, x' \rangle)^k$.

תכונות של תת-גרייאנט

סאב-גרייאנט: וקטור u הוא סאב-גרייאנט של f אם: $\forall u, f(u) \geq f(w) + \langle u, w - u \rangle$

סימון: $\partial f(w)$ הוא קבוצת כל סאב-גרייאנטים של f בנקודה w .

למה: f קמורה אם"ם לכל נקודה $w \in \text{dom}(f)$ מתקיים $\partial f(w) = \{ \nabla f(x) \mid x \in \text{epi}(f) \}$ אם f דפרנציבילית ב- w אז $\partial f(w) \neq \emptyset$.

טענה: ניח $f: V \rightarrow \mathbb{R}$ קמורה, ותהי $f(x) = \max_{j \in J} f_j(x)$. כמו כן עבור $R \subseteq \partial f(w)$ ו- $J \in \arg \max_{j \in J} f_j(w)$ עבור $w \in V$ הו"ע של $\partial f(w)$, $\vec{0} \in \partial f(w)$ לכל $w \in V$ הו"ע של $\partial f(w)$.

ליפשיץ: פונקציה $f: C \rightarrow \mathbb{R}$ נקראת ρ -ליפשיץ אם: $\forall w_1, w_2 \in C \quad |f(w_1) - f(w_2)| \leq \rho \|w_1 - w_2\|$

למה: אם f קמורה, אז f היא ρ -ליפשיץ אם"ם הנורמה של כל סאב-גרייאנט של f הוא לכל היותר ρ .

אופטימיזציה קמורה

אופטימיזציה: $\min_x f_0(x)$ s.t. $f_i(x) \leq b_i \quad \forall i = 1, \dots, m$

בעיית אופטימיזציה קמורה: בעיית אופטימיזציה f_i קמורה. בעיית תכנון לינאר: בעיית אופטימיזציה f_i קמורה ב- \mathcal{H} . דוגמה: $ERM_{M_{\mathcal{H}}}$.

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

טענה: לא כל הבעיות הקמורות מעל \mathbb{R}^d הן למידות-PAC.

טענה: בעיית קמורה מעל \mathbb{R}^d ו- ℓ סימטרי ליפשיץ היא למידה-PAC.

- \mathcal{H} **סימטרי** - קיים B עבור $\forall w \in \mathcal{H} \quad \|w\| \leq B$
- \mathcal{H} **ליפשיץ** - פונקציה $\ell(x, y)$ קמורה ב- \mathcal{H} הוא ρ -ליפשיץ.

למקרה זה סיבוכיות הדגימה תלויה רק ב- ρ, B, ϵ .

Sub-gradient Descent Algorithm:

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

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

- Choose $v_t \in \partial f(w^{(t)})$
- $w^{(t+1)} \leftarrow w^{(t)} - \eta \cdot v_t$

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

משפט: תהי f פונקציה קמורה, ρ -ליפשיץ, ויהי $w^* \in \arg \min_{w \in \mathcal{H}} f(w)$. אם נריץ את האלגוריתם SGD על f במשך T צעדים $\frac{\|w^*\|^2}{\rho^2 T}$, η אז הקוטור \bar{w} המוחזר מקיים:

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

מסקנה: תהי f פונקציה קמורה, ρ -ליפשיץ, ויהי $w^* \in \arg \min_{w \in \mathcal{H}} f(w)$. אם $\epsilon > 0$ אם נריץ את האלגוריתם SGD למעור f במשך $\frac{\|w^*\|^2 \rho^2}{\epsilon^2}$ צעדים $\frac{\|w^*\| \rho}{\sqrt{T}}$, אז הקוטור \bar{w} המוחזר מקיים: $f(\bar{w}) \leq f(w^*) + \epsilon$

מסקנה: Sub-GD צורך $\frac{\|w^*\|^2 \rho^2}{\epsilon^2}$ איטרציות בכדי להתכנס.

Stochastic Gradient Descent Algorithm:

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

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

- Choose $(x, y) \sim \mathcal{D}$
- Choose $v_t \in \partial \ell(w^{(t)}(x, y))$
- $w^{(t+1)} \leftarrow w^{(t)} - \eta \cdot v_t$

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

למה: תהי בעיית למידה קמורה-ליפשיץ-סימטריה עם פרמטרים p, B . לכל $\epsilon > 0$, אם נריץ את האלגוריתם SGD למעור $L_D(w)$ עם מספר איטרציות $T \geq \frac{B^2 \rho^2}{\epsilon^2}$ ו- $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$, אז הפלט SGD מקיים:

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

טענה: אלוגוריתם \mathcal{A} תלומד בעזרת SGD הוא לומד-חלש.

פונקציות אקטיביות מוכרות:

Sigmoid(x) $= \frac{1}{1+e^{-x}}$, **Tanh(x)** $= \frac{2}{1+e^{-2x}} - 1$

ReLU(x) $= \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$

באלגוריתם SGD רשתות בוחרים את w ההתחלית באופן **רנדומלי** (עם התפלגות כך ש- $w^{(1)}$ קרוב מספיק ל-0).

Backpropagation Algorithm:

input: example (x, y) , weight vector W , layered graph (V, E) and 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 $W_{t,l,j}$ as weight of $(v_{t,l}, v_{t+1,j})$

forward: set $O_0 = x$ and for $t = 1, \dots, T$:

$$\text{set } a_{t,l} = \sum_{i=1}^{k_{t-1}} W_{t-1,l,i} \cdot O_{t-1,i} + O_{t-1,l} = \sigma(a_{t,l})$$