

- ערכים גבוהים או נמוכים מדי יגרמו לפונקציה לעבוד בתחום לא מתאים (גרדיאנט 0)

#### שיטות נורמליזציה:

- Min-Max:  $v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new\_max}_A - \text{new\_min}_A) + \text{new\_min}_A$
- Z-Score:  $v' = \frac{v - \text{mean}_A}{\text{stand\_dev}_A}$

#### אלגוריתמים לאופטימיזציות:

נרצה לעשות את הלמידה בצעדים קטנים יותר (ללא קפיצות גדולות). נשתמש בקצב למידה (LR).

- GD
- SGD
- Mini-Batch GD
- Momentum
- AdaGrad
- RMSprop
- Adam

#### :GD

מטרה: ביצוע אופטימיזציה. מזעור ה-Loss.

יש  $w$  וקטור משקולות שבעזרתו ניתן משקל לכל פיצ'ר.

#### תכונות:

- קירוב טוב לגרדיאנט
- צריך לעבור בכל איטרציה על כל הדוגמאות
- מתכנס בצורה איטית

$$w^t \leftarrow w^{t-1} - \text{eta} * \frac{\partial L}{\partial w^{t-1}} \quad \text{where } w^t \text{ is the vector in iteration } t$$

בכל עדכון מחשבים את הגרדיאנט על פי הממוצע על כל דוגמאות סט האימון.

ההתכנסות איטית (לאפוק).

#### :SGD

#### תכונות:

- נדגום רנדומלית דוגמא אחת בכל איטרציה
- לאחר מספיק איטרציות האלגוריתם יתכנס

אותה פונקציית עדכון.

## :GD VS SGD

GD איטי יותר. SGD יותר "רועש" ופחות צפוי.

## :Mini-Batch GD

דומה ל-SGD אך בכל איטרציה נדגום k דוגמאות אקראיות וב"ת.

## :Momentum

ל-SGD קשה באיזורים בהם המשטח של הדוגמאות מתעקל בצורה תלולה יותר (עם שיפוע תלול יותר), במימד אחד יותר מאשר באחר. במקרים האלו SGD עובר במישור בתנודות גדולות אך עם התקדמות מאוד איטית לאופטימום המקומי. לכן נרצה להשתמש ב-momentum.

זוהי בעצם מתודה שהיא תוספת ל-SGD שעוזרת לכוון את וקטורי הגרדיאנטים לכיוון הנכון מהר יותר (התכנסות מהירה יותר). שיטה זו כמעט תמיד עובדת טוב יותר ומהר יותר מאשר SGD.

כלל עדכון:

$$v_t = \gamma v_{t-1} + \eta * \frac{\partial L}{\partial \theta_i}, \quad \theta_i = \theta_i - v_t$$

## :AdaGrad

אלגוריתם לאופטימיזציה מבוססת גרדיאנטים שמתאים את קצב הלמידה (eta) לכל פרמטר.

- נפחית את קצב הלמידה של משקולות עם גרדיאנטים גבוהים.
- נעלה את קצב הלמידה של משקולות שמקבלות עדכונים קטנים או לא מתעדכנים בצורה עקבית.

כלל עדכון:

$$c_{t,i} = \left( \sum_{k=0}^t \frac{\partial L}{\partial \theta_{k,i}} \right)^2, \quad \theta_{t+1,i} = \theta_{t,i} - \frac{\frac{\partial L}{\partial \theta_{t,i}}}{\sqrt{c_{t,i} + \epsilon}}$$

## :RMSprop

במקרה של deep learning ל-AdaGrad יש מגרעה, השינויים בקצב למידה שלו לרוב הם מאוד אגרסיביים וגורמים לעצירת הלמידה מוקדם מדי. לכן במקרה הזה נשתמש בשיטה הזו.

כלל עדכון:

$$c_{t,i} = \gamma c_{t-1,i} + (1 - \gamma) \left( \frac{\partial L}{\partial \theta_{t,i}} \right)^2, \quad \theta_{t+1,i} = \theta_{t,i} - \frac{\frac{\partial L}{\partial \theta_{t,i}}}{\sqrt{c_{t,i} + \epsilon}}$$

ערכים טיפוסיים ל- $\gamma$  הם 0.9, 0.99, 0.999.

שלא כמו AdaGrad עדכונים לא נהיים קטנים יותר מנוטוניות.

## Adam

בעקרון זה שילוב של momentum ו-RMSprop. מעשית זה האלגוריתם הדיפולטיבי בהרבה חבילות.

כלל עדכון:

$$m_{t,i} = \gamma_1 m_{t-1,i} + (1 - \gamma_1) \frac{\partial L}{\partial \theta_{t,i}}$$

$$c_{t,i} = \gamma_2 v_{t-1,i} + (1 - \gamma_2) \left( \frac{\partial L}{\partial \theta_{t,i}} \right)^2$$

$$\theta_{t+1,i} = \theta_{t,i} - \frac{m_{t,i}}{\sqrt{v_{t,i} + \epsilon}}$$

ערכים מומלצים הם:  $\gamma_1 = 0.9, \gamma_2 = 0.99$

### למידת PAC (PAC learnability):

נגדיר  $\delta$  – הסיכוי שנעבוד על סט דוגמאות לא מייצג.

$\epsilon$  – אי דיוק הסיווג.

$$E_{x,y}[l(y, f_\theta(x))] > \epsilon \Rightarrow \text{classifier failed}$$

$$E_{x,y}[l(y, f_\theta(x))] \leq \epsilon \Rightarrow \text{classifier approximately correct}$$

סט היפותיזות  $H$  נקרא למידת PAC ביחס לפונקציית  $\bar{l}(x, y; h)$  אם קיים מס' מינימלי של דוגמאות אימון  $m$  שמוגרל בצורת  $B^n$  ומאותה התפלגות מהסתברות  $D$  ומתקיים:

$$(x, y) \sim D, E[\bar{l}(x, y; h)] \leq \min_{h' \in H} E[\bar{l}(x, y; h')] + \epsilon$$

זה אומר שקיים סט היפותיזות שעבורו אם נקבל מספיק דוגמאות ( $m$ ) אז נצליח למצוא סיווג טוב בשגיאה עד  $\epsilon$  בהסתברות  $1 - \delta$ .

משפט: סט היפותיזות סופי הוא למידת PAC.

$$P_S \left( \left| \frac{1}{m} \sum_{i=1}^m \bar{l}(x_i, y_i; h) - E[\bar{l}(x, y; h)] \right| < \epsilon \right) > 1 - \delta$$

במילים: דיוק גבוה של הלמידה  $\leftarrow$  ההפרש בין הממוצע על סט האימון לבין תוחלת הדוגמאות שלא ראינו הוא קטן.

### סוגי אלגוריתמים:

- KNN
- K Means
- Perceptron
- bayesian network
- SVM
- Convolutional neural network (CNN)