

# תרגיל 4 - חסדק תיאורטי

על

1.1

1) יהי  $\alpha \in [0, 1]$

$$\alpha g(u) + (1-\alpha) g(v) = \alpha \sum_{i=1}^m \tau_i f_i(u) + (1-\alpha) \sum_{i=1}^m \tau_i f_i(v)$$

הגזיות  $g$

$$= \sum_{i=1}^m [\alpha \tau_i f_i(u) + (1-\alpha) \tau_i f_i(v)]$$

$$\geq \sum_{i=1}^m \tau_i f_i(\alpha u + (1-\alpha)v) = g(\alpha u + (1-\alpha)v)$$

הגזיות  $g$

$$\alpha f_i(u) + (1-\alpha) f_i(v) \geq f_i(\alpha u + (1-\alpha)v)$$

כל  $i \in \{1, \dots, m\}$  מתקיים מהנחות  $\alpha, 1-\alpha \in [0, 1]$  ו- $\tau_i > 0$  וכן

$$\tau_i (\alpha f_i(u) + (1-\alpha) f_i(v)) \geq \tau_i f_i(\alpha u + (1-\alpha)v)$$

$$F(x) = -x, g(x) = x^2 \quad \text{2) נקודות קיצון}$$

$$F \circ g(x) = -x^2 \text{ ו-} F \circ g(x) = -x^2 \text{ איננו המורה (אם לא קצרה)}$$

$$F = g = e^{-x} \text{ ו-} F = g = e^{-x} \text{ איננו המורה (אם לא קצרה)}$$

1.2

3) המושג שטען להקדירה היא המקסימום בין שתי פונקציות שלישיות  $w, w$  בפני קמורות

וכפי שהאמינו בכיתה (בשקופית 4) בגרסא  $g$  של  $f$  (זכור  $f$ ), המקסימום של  $g$  הוא פונקציה קמורה

היא פונקציה קמורה בעצמה, בפרט  $f$  שהגזית

4) בתמונה שהבטנו נכון שהנוחה היא אך ורק להגדיר פונקציה  $g$   $g(x, y) = 0$   $g(x, y) = 0$   $g(x, y) = 0$

$$g = \begin{pmatrix} 0 & , l_{x,y}(w,b) = 0 \\ (-y, -y) & , l_{x,y}(w,b) = 0 \end{pmatrix}$$

$$F_K(y) \geq F_K(x) + \langle g_K(x), y-x \rangle = (g_K \in \partial F_K(x)) \text{ sub-gradient}$$

$$F(y) = \sum_{i=1}^n F_K(y) \geq \sum_{i=1}^n [F_K(x) + \langle g_K(x), y-x \rangle] = \sum_{i=1}^n F_K(x) + \sum_{i=1}^n \langle g_K(x), y-x \rangle$$

$$= F(x) + \langle \sum_{i=1}^n g_K(x), y-x \rangle = F(x) + \langle g(x), y-x \rangle$$

$$\sum_{i=1}^n g_K(x) = g(x) \in \partial F(x) = \partial \sum_{i=1}^n F_K(x) \text{ - sub-gradient}$$

⑥ הוכחה ש- $g_i$  הוא sub-gradient של  $\ell_{x_i, y_i}^{hinge}(w, b)$  בנקודה  $(w, b)$  שבה  $\ell_{x_i, y_i}^{hinge}(w, b) > 0$

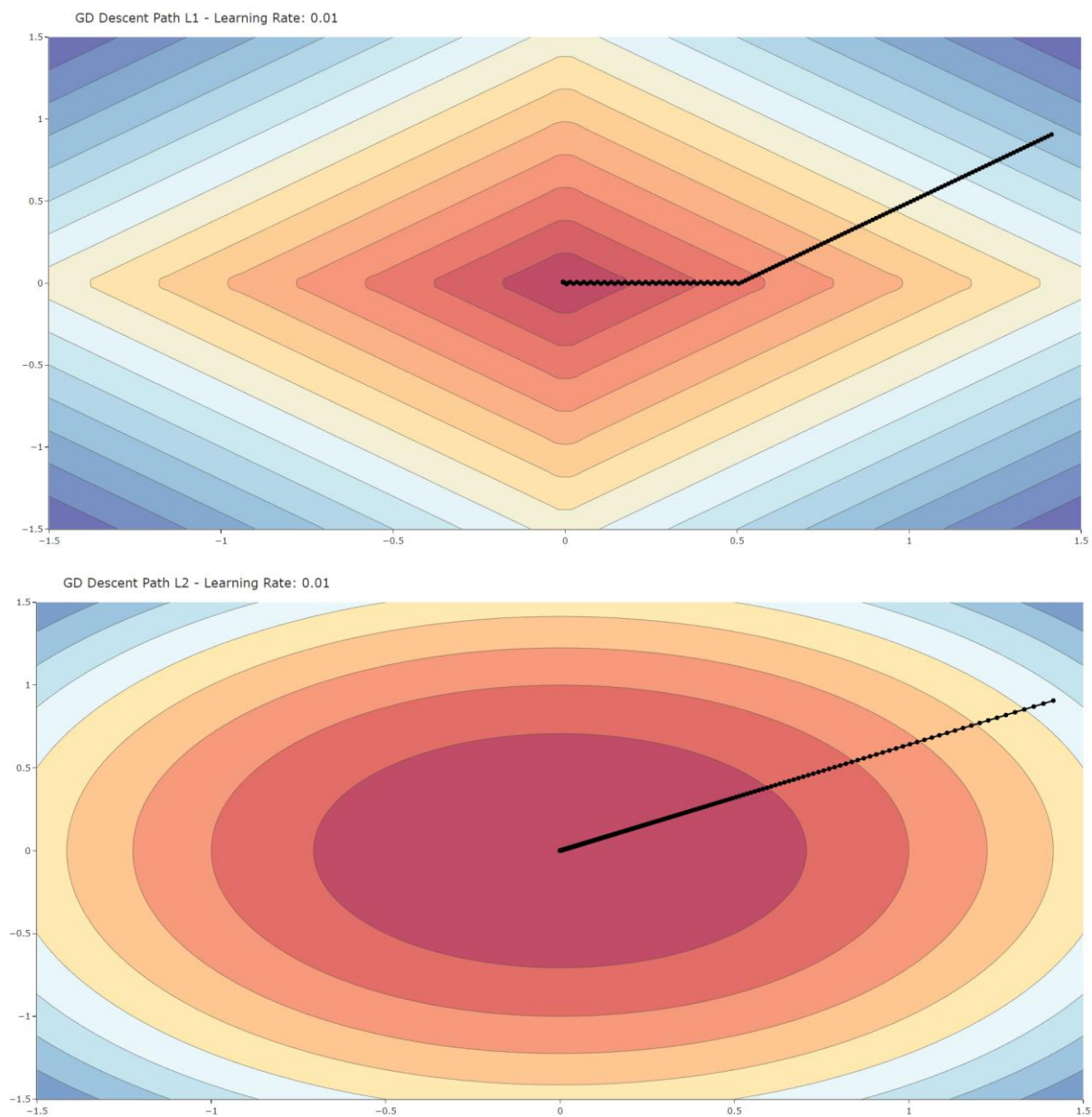
הוכחה: נניח  $\ell_{x_i, y_i}^{hinge}(w, b) > 0$

$$\frac{1}{m} \sum_i g_i + \alpha(w, b) \frac{\lambda \|w\|^2}{2} \leq \partial \frac{1}{m} \sum_i \ell_{x_i, y_i}^{hinge}(w, b) + \frac{\lambda}{2} \|w\|^2$$

## חלק פרקטי

### 2.1

.1

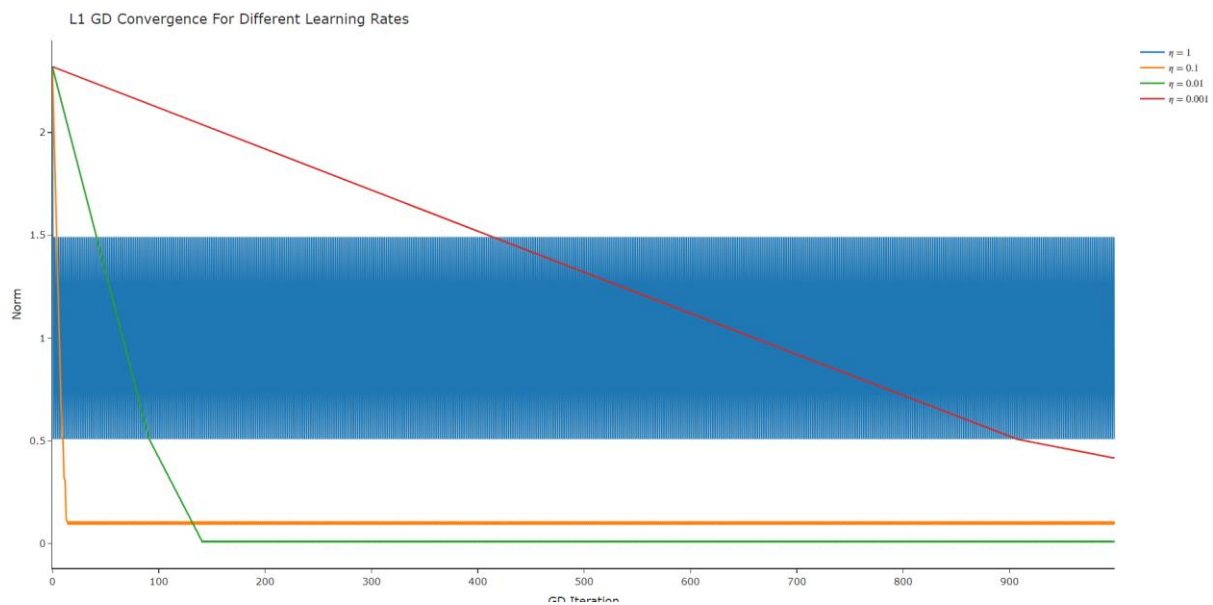


ההבדל המרכזי שניתן לזהות הוא שעבור  $L_2$  נקבל תנועה אחידה עבור כל רמה

2. נשים לב ל-2 תופעות:

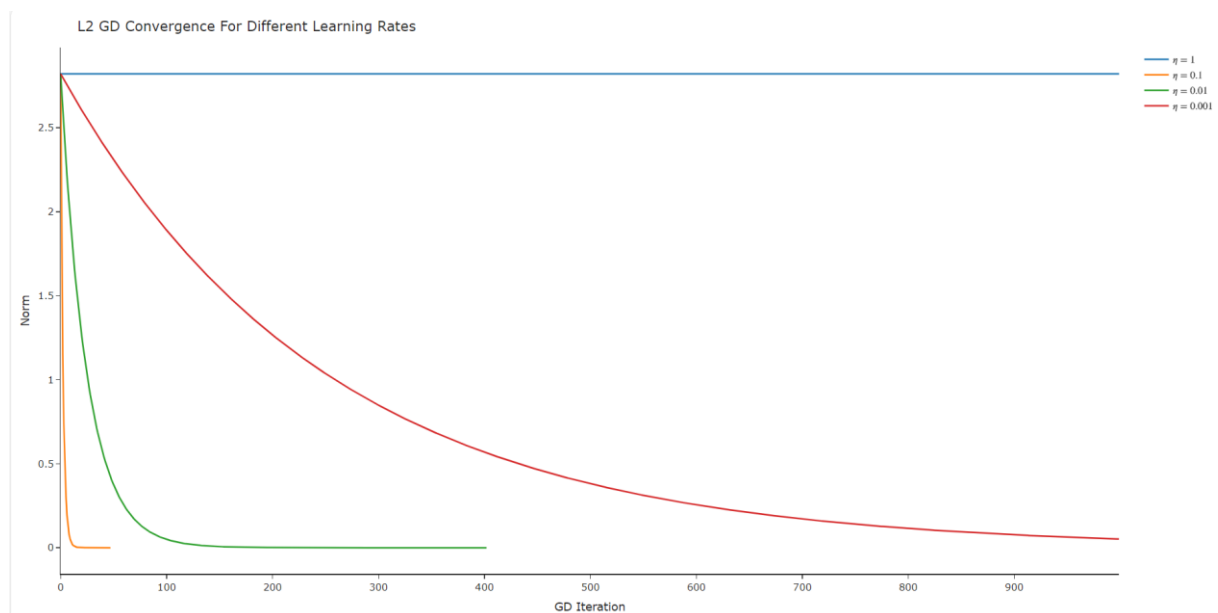
- שינוי הכיוון של הירידה מירידה אלכסונית להמשך על הציר המקביל.
- לאחר ההמשך בציר במקביל יש קפיצה בין 2 ערכים – כלומר, כל פעם מפספסים את ההתכנסות ומנסים לחזור, אך שינוי הכיוון הוא חד מדי. המשמעות היא שבפועל הגענו למקסימום ההתכנסות שנוכל.

3.



עבור  $\eta=1$  אנחנו מקבלים over-shoot חוזר של המודל, שכל פעם קופץ בצורה קיצונית מעל המינימום של פונקציית ה-Loss מה שמוביל לאוסילציות אין סופיות.

מעבר לכך אנחנו רואים שעבור  $\eta=0.1$  ההתכנסות היא מאוד מהירה, ואילו ככל שהוא קטן ההתכנסות היא הרבה יותר איטית (נקווה שב-0.001 הוא אכן מתכנס בסוף 😊...)



עבור  $\eta=1$  קבלנו התכנסות לפתרון לא טוב (נורמה מאוד גבוהה שלא השתפרה כלל). לשאר האטות נקבל אותו דבר כמו קודם, רק שכאן הלמידה הרבה יותר מהירה.

4.

```
for learning rate: 1 - Lowest loss achieved for L1 module: 0.5081196195534134
for learning rate: 0.1 - Lowest loss achieved for L1 module: 0.09188038044658689
for learning rate: 0.01 - Lowest loss achieved for L1 module: 0.008119619553413011
for learning rate: 0.001 - Lowest loss achieved for L1 module: 0.41611961955345933
for learning rate: 1 - Lowest loss achieved for L2 module: 2.8210062332145176
for learning rate: 0.1 - Lowest loss achieved for L2 module: 2.19211242161629e-09
for learning rate: 0.01 - Lowest loss achieved for L2 module: 2.4898254541043804e-07
for learning rate: 0.001 - Lowest loss achieved for L2 module: 0.05166846244106496
```

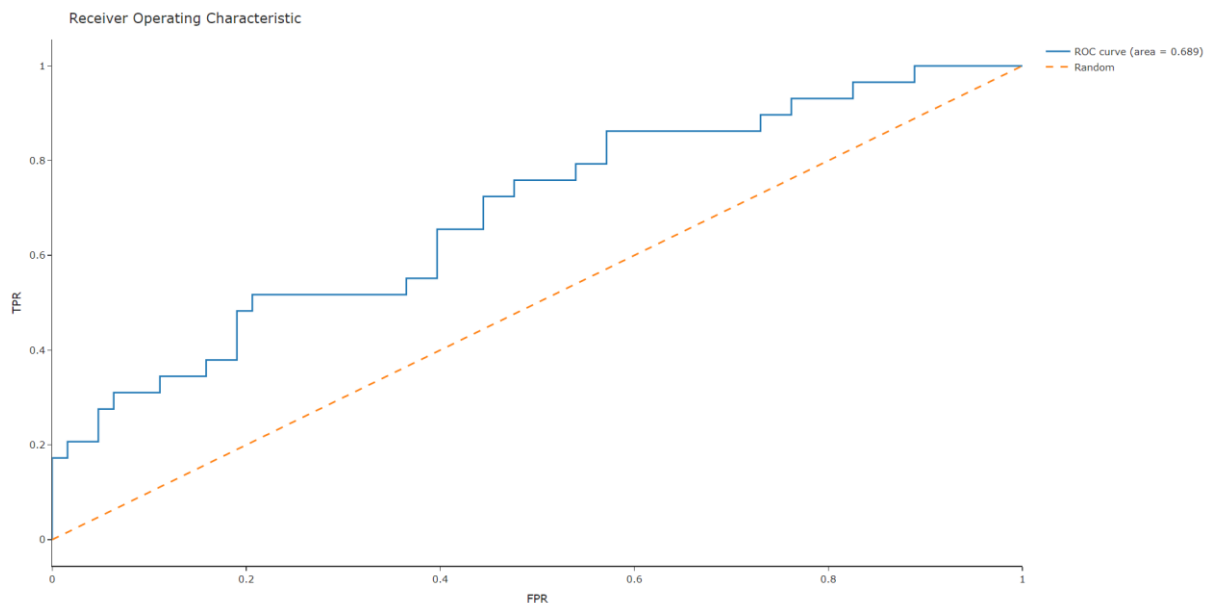
למה L1 פחות טוב מL2?

קודם כל, L1 מנחית חלק מהמשקולים ל0, מה שגורם לפתרון יותר דליל שמאבד חלק מהיכולת להסיק דברים בעזרת פיצ'רים מסויימים (בפרט במקרים בו יש קשר ביניהם).

אך הסיבה העיקרית היא שהערך שהיא מאפטמת את החציון כלומר גרדיאנט קבוע (הנגזרת של פולינום ממעלה ראשונה הוא קבוע) ואילו L2 מאפטמת את הממוצע (משתנה, כיוון שהיא תלויה בx) ולכן הנורמה משתנה כתלות ב X ולכן למודל יש יכולת ללמוד ולהתאים את עצמו ל loss קטן יותר ולהתקרב יותר להתכנסות.

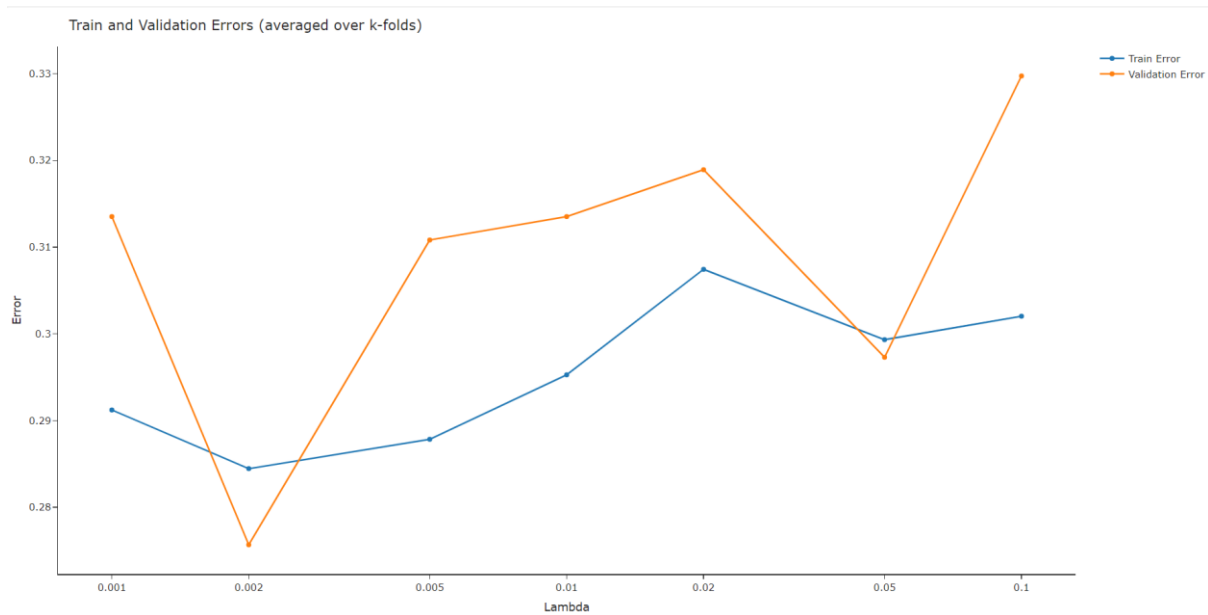
## 2.2

.5



.6

```
Optimal alpha: 0.4785254786906066
Test error with optimal alpha: 0.30434782608695654
```



Optimal Regularization Parameter: 0.002  
Model achieved test error of 0.25