

Theoretical Part

1.1 Convex optimization

1. נתון f_1, \dots, f_m פונקציות קמורות כלליות $\alpha \in [0, 1]$

כך שכל $i \in [m]$! $C \ni y, x$

$$f_i(\alpha x + (1-\alpha)y) \leq \alpha f_i(x) + (1-\alpha)f_i(y)$$

כעת נראה זאת על g קמורה. $\bar{\alpha} \in [0, 1]$! $C \ni y, x$

$$g(\alpha x + (1-\alpha)y) = \sum_{i=1}^m \theta_i f_i(\alpha x + (1-\alpha)y)$$

$$\leq \sum_{i=1}^m \theta_i (\alpha f_i(x) + (1-\alpha)f_i(y))$$

$$= \sum_{i=1}^m \theta_i \alpha f_i(x) + \sum_{i=1}^m \theta_i (1-\alpha) f_i(y)$$

$$= \alpha \sum_{i=1}^m \theta_i f_i(x) + (1-\alpha) \sum_{i=1}^m \theta_i f_i(y)$$

$$= \alpha g(x) + (1-\alpha)g(y)$$

כנראה.

2. דוגמה נכזית: נגדיר $g(x) = x^2$, $f(x) = -x$ שתי הפונקציות

הן קמורות ע"י האזנה אפס מתק"פ - $h = f \circ g = -(x^2)$

אין קמורה.

1.2 sub-gradients for soft - SVM objective

3. ראינו כייתה $e - b + w^T x$ (פונקציה אפנית) היא פונקציה קמורה ואכן נקרא $e - (b + w^T x)$ היא גם קמורה. 0 היא פונקציה קמורה באופן טריוויאלי. וראינו כייתה כי פונקציית \max של פונקציות קמורות היא קמורה גם כן. אכן סה"כ נקרא $e - \text{hinge loss}$ היא פונקציה קמורה.

4. בהינתן $\ell_{\text{hinge}}(w, b)_{x,y}$ מתקיים:

$$\ell_{\text{hinge}}(w, b)_{x,y} = \begin{cases} 0 & y(x^T w + b) \geq 1 \\ 1 - y(x^T w + b) & y(x^T w + b) < 1 \end{cases}$$

ה- $\text{sub-gradient } g$ מוגדר כאופן הקרא: $g(\frac{\partial \ell_{\text{hinge}}}{\partial w}, \frac{\partial \ell_{\text{hinge}}}{\partial b})$
כל נפסד מתקיים: $y(x^T w + b) \geq 1 \leftarrow \ell_{\text{hinge}}(w, b) = 0$ כל:

$$\frac{\partial \ell_{\text{hinge}}}{\partial w} = 0, \frac{\partial \ell_{\text{hinge}}}{\partial b} = 0 \Rightarrow g = 0$$

כל $1 - y(x^T w + b) = \ell_{\text{hinge}}(w, b) \leftarrow y(w^T x + b) < 1$ כל:

$$\frac{\partial \ell_{\text{hinge}}}{\partial w} = -yx, \frac{\partial \ell_{\text{hinge}}}{\partial b} = -y \Rightarrow g(-yx, -y)$$

$$g = \begin{cases} 0 & \text{if } \ell_{\text{hinge}}(w, b) = 0 \\ (-yx, -y) & \text{if } \ell_{\text{hinge}}(w, b) = 1 - y(x^T w + b) \end{cases}$$

5. נחשברה: $\partial f_k(x) \ni g_k$ כל $k \in [m]$ כל $y \in \mathbb{R}^d$:
 $f_k(y) \geq f_k(x) + g_k^T(x - y)$ פ"ק נ"כ

$$\sum_{i=1}^m f_i(y) \geq \sum_{i=1}^m f_i(x) + g_i^T(x - y) = \sum_{i=1}^m f_i(x) + \sum_{i=1}^m g_i^T(x - y)$$

נחשברה: $f(x) = \sum_{i=1}^m f_i(x)$ פ"ק נ"כ

$$f(y) \geq f(x) + \sum_{i=1}^m g_i^T(x - y)$$

$$\rightarrow f(y) \geq f(x) + \left(\sum_{i=1}^m g_i \right)^T (x - y)$$

כנראה.

6. נשים לב כי מלפני 4 חשבו וקראנו:

$$g_i = \begin{cases} 0 & \text{if } \ell_{\text{hinge}}(w, b) = 0 \\ (-y_i x_i, -y_i) & \text{else} \end{cases}$$

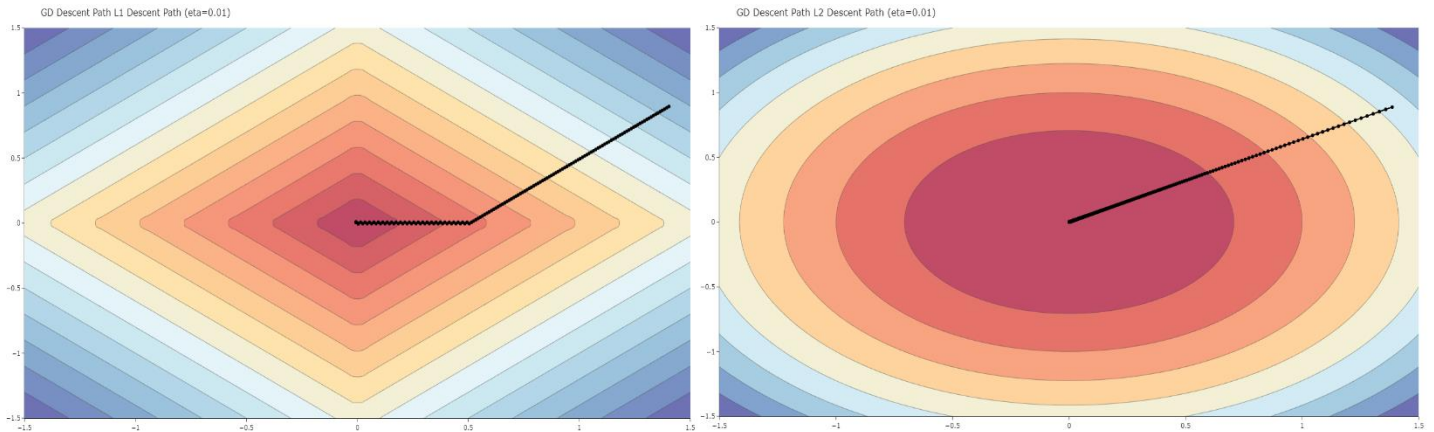
בנוסף, מלפני 5, שתי הפונקציות סכום הן קמורות. וכפ"ל סקרא
 על פונקציה הו"א קמור נקרא:

$$\frac{1}{m} \sum_{i=1}^m g_i + \lambda \|w\| \in \partial f(w, b) \quad w \in \mathbb{R}^d$$

Practical part – IML

2.1.1 Comparing Fixed learning rates

1. explain the differences seen between the L1 and L2 modules:



נשים לב להבדלים בין המודלים:

במודל L2 הכיוון של הגרדיאנט הוא בקו ישר אל נקודת המינימום בעוד שבמודל L1 הגרדיאנט יורד בקו ישר עד נקודה כלשהי (עד שערך ה γ מתאפס) ואז מבצע פנייה חדה ויורד בזיג-זג במקביל לציר ה-X עד נקודת המינימום בערך $(0,0)$.

הדבר נובע מכך שבמודל L2 כל שתי נקודות שהן על אותו "טבעת" יש מרחק שווה מנקודת המינימום לכן ההתנהגות של הקו הוא לרדת במורד הטבעות ישירות אל נקודת המינימום בעוד שבמודל L1 לכל שתי נקודות שנמצאות על אותו "מסגרת" אין מרחק שווה מנקודת המינימום.

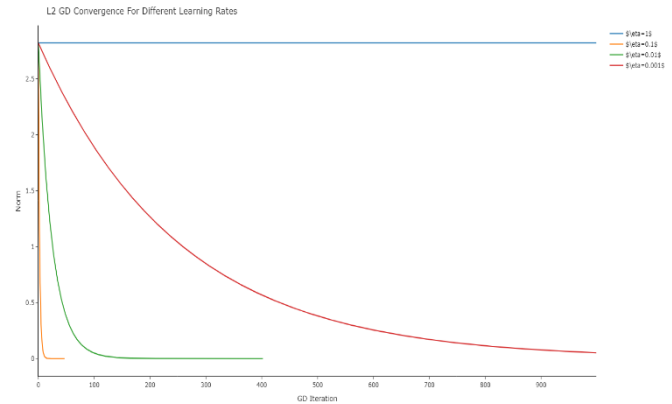
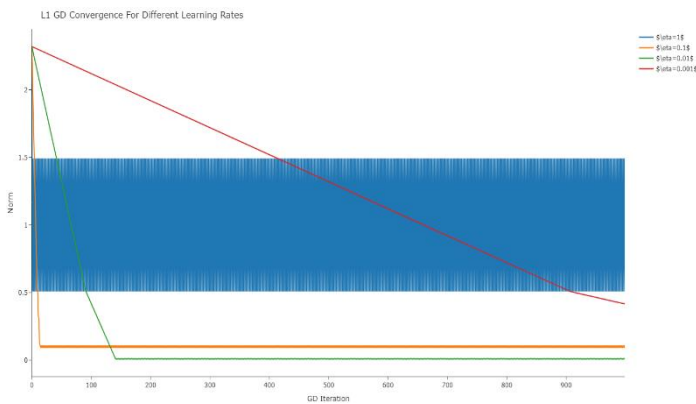
2. following the previous question describe two phenomena that you have seen in the descent path of the ℓ_1 objective when using GD and a fixed learning rate:

שתי תופעות שמתרחשות במודל L1 הן:

1. כיוון הגרדיאנט יורד באלכסון לכיוון נקודת המינימום עד שהוא מבצע "שבירה" כאשר ערך ה γ מתאפס ואז הכיוון הוא לאורך ציר ה אעד נקודת המינימום.

2. מטבע התנהגות הפונקציה L1, כאשר מתקרבים אל נקודת המינימום, הפונקציה מזגזגת אל המינימום (עולה ויורדת) ולא יורדת בקו ישר.

3. For each of the modules, plot the convergence rate (i.e. the norm as a function of the GD iteration) for all specified learning rates. Explain your results:



נסתכל על הגרף של המודל L1 (צד שמאל)

$\eta = 1$: הנורמה נותרת תנודתית, וערכי הגרדיאנט נעים ללא התכנסות מ-0.5 ל-1.5. מה שמעיד על קצב למידה גבוה מדי שגורם לקפיצות יתר ולאי התכנסות.

$\eta = 0.1$: הנורמה יורדת במהירות לנקודה יציבה בתוך מספר קטן של איטרציות ונשארת קבועה, מה שמעיד על קצב למידה יחסית טוב.

$\eta = 0.01$: הנורמה יורדת בהדרגה ומתכנסת לנקודה יציבה לאחר כ-140 איטרציות לערך.

$\eta = 0.001$: הנורמה יורדת מאוד לאט לאורך כל האיטרציות, מה שמעיד על קצב התכנסות איטי.

נסתכל על הגרף של המודל L2 (צד ימין)

$\eta = 1$: הנורמה לא יורדת, מה שמעיד על קצב למידה גבוה מדי וכישלון בהתכנסות.

$\eta = 0.1$: הנורמה יורדת במהירות לאפס בתוך 50 איטרציות, מה שמראה על התכנסות מהירה מדי.

$\eta = 0.01$: הנורמה יורדת בהדרגה ומגיעה כמעט לאפס לאחר כ-200 איטרציות בערך.

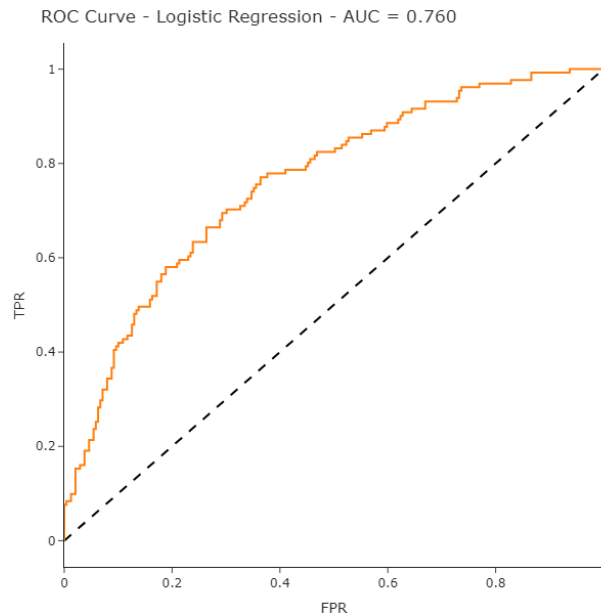
$\eta = 0.001$: הנורמה יורדת בהדרגה לאורך כל האיטרציות, מה שמעיד על קצב התכנסות איטי מדי.

נסיק מהגרפים שקצב הלמידה אופטימלי חשוב להתכנסות יעילה. קצב איטי מדי עלול לקחת זמן רב וקצב מהיר מדי עלול לגרום ל"קפיצות". קצב הלמידה משתנה בין שני המודלים L1 ל-L2 בשל השפעתם השונה על עדכוני המשקולות במהלך גרדיאנט דסנט.

4. What is the lowest loss achieved when minimizing each of the modules? Explain the differences:

2.2 Minimizing Regularized Logistic Regression

5. Using your implementation, fit a logistic regression model over the data. Use the `predict_proba` to plot an ROC curve:



6. Which value of α achieves the optimal ROC value according to the criterion below. Using this value of α * what is the model's test error?:

```
pnina_ei@pond:~/IML/ex4 $ python3 gradient_descent_investigation.py
The value of  $\alpha$  achieves the optimal ROC is: 0.3247731561147603
The model's test error is: 0.33695652173913043
```

קבלנו שהאלפא האופטימלית היא: $\alpha = 0.324$ ושגיאת הטסט היא: **0.336**.

7. What value of λ was selected and what is the model's test error?:

```
Optimal regularization parameter: 0.02
Model achieved test error of 0.28
```

קבלנו שהלמדא האופטימלית היא $\lambda = 0.02$. ושגיאת הטסט היא **0.28**.