

1.1.1

$\alpha \in [0,1]$, $v, w \in C$ הינה, ורנו $R - S$ כ- N המבוקש מינ' ג' ' נסוכ

$$g(\alpha v + (1-\alpha)w) = \sum_{i=1}^m \gamma_i f_i(\alpha v + (1-\alpha)w) \leq \sum_{i=1}^m \gamma_i (\alpha f_i(v) + (1-\alpha) f_i(w))$$

↓
 ↗
 f_i

$$= \alpha \sum_{i=1}^m \gamma_i f_i(v) + (1-\alpha) \sum_{i=1}^m \gamma_i f_i(w) = \alpha g(v) + (1-\alpha) g(w) \Rightarrow \boxed{\text{הינו ג' } C \text{-נסוכ}}$$

1.1.2

$$f(x) = |x|, g(x) = x^2 - 1$$

$\left\{ \begin{array}{l} \text{נורמליזציה} \\ \text{הינה } "1" \text{ נורמליז.} \\ \text{הנה } h(x) = |x^2 - 1| \end{array} \right.$

$$\Rightarrow h(x) = |x^2 - 1|$$

$R \subset \mathbb{C}$ ו $v, w \in R$, $\alpha \in [0,1]$ "

כדאי ב' "



$$f(\alpha v + (1-\alpha)w) = |\alpha v + (1-\alpha)w| \leq \alpha |v| + (1-\alpha) |w| = \alpha f(v) + (1-\alpha) f(w) \Rightarrow \underline{f}$$

$$\text{לדוגמא } 2-3 \text{ ו } 1+2 \quad g(x) = x^2 > 0 \quad \Rightarrow \text{ והם יוצרים}$$

הינה g , \underline{g}

$$\alpha = \frac{1}{2} \quad \therefore \quad v, w = -1, 1 \quad \text{ל' } \underline{g}$$

$$h(w) = |1-1| = 0, h(v) = |1+1| = 2, (1-\alpha)v = \frac{1}{2}, \alpha \cdot v = -\frac{1}{2} \Leftarrow$$

$$h(\alpha v + (1-\alpha)w) = 2(-\frac{1}{2} + \frac{1}{2}) = h(0) = \underline{1}, \alpha h(v) + (1-\alpha)h(w) = \underline{0} \Leftarrow$$

$$h(\alpha v + (1-\alpha)w) = 1 > 0 = \alpha h(v) + (1-\alpha)h(w) \Leftarrow$$

- יונט איזייד $\underline{\underline{h}} \Leftarrow$

1.2.3

הנחתה דענו כי $\max_{\alpha \in [0,1]} f(\alpha v + (1-\alpha)w)$ מינימום של פונקציית האיחוד $f+g$: $C \rightarrow \mathbb{R}$

$$\begin{aligned} & \max_{\alpha \in [0,1]} \{f(\alpha v + (1-\alpha)w), g(\alpha v + (1-\alpha)w)\} \leq \max \{f(v) + (1-\alpha)f(w), g(v) + (1-\alpha)g(w)\} \\ & \leq \alpha \cdot \max \{f(v), g(v)\} + (1-\alpha) \max \{f(w), g(w)\} \end{aligned}$$

הנחתה $\exists \alpha \in [0,1]$ בפונקציית האיחוד $f+g$ מינימום של $1-y(x^T w + b)$ נובעת מכך ש f ו- g הן פונקציות יריריות.

$$f(\alpha v + (1-\alpha)w) = A(\alpha v + (1-\alpha)w) + b \leq \alpha Av + b + (1-\alpha)Aw + b = \alpha f(v) + (1-\alpha)f(w)$$

הנחתה $\exists \alpha \in [0,1]$ מינימום של $1-y(x^T w + b)$ מינימום של $1-y(x^T w + b)$.

1.2.4

$$\frac{\partial}{\partial w} (1-y(x^T w + b)) = -y_x, \quad \frac{\partial}{\partial b} (1-y(x^T w + b)) = -y$$

$$\lambda(w, b) = \max \{0, 1-y(x^T w + b)\}$$

$$g = (0, 0) \text{ מינימום של } \lambda \text{ בסביבת } h(w, b) = 0 \Leftrightarrow 1-y(x^T w + b) < 0 \quad (1)$$

$$g = (-y_x, -y) \text{ מינימום של } \lambda(w, b) \text{ בסביבת } h(w, b) \Leftrightarrow 1-y(x^T w + b) > 0 \quad (2)$$

$$\text{במקרה השני } h \leq 1-y(x^T w + b) = 0 \quad (3)$$

הנחתה $\exists \alpha \in [0,1]$ מינימום של λ בסביבת $h(w, b) = 0$ מינימום של $1-y(x^T w + b)$ בסביבת $h(w, b) = 0$.

$$\lambda \in [0,1] \quad g = (-\lambda y_x, -\lambda y)$$

$$g = \begin{cases} (0) & , 1-y(x^T w + b) < 0 \\ (-y_x, -y) & , 1-y(x^T w + b) \geq 0 \\ (-\lambda y_x, -\lambda y) & , 1-y(x^T w + b) = 0 \end{cases}$$

1.2.5

לעתה $f_i(x) \in \mathbb{R}^m$ if $i \in [n]$ סביר $x, w \in \mathbb{R}^d$

$$\forall i \in [n] \quad f_i(u) \geq f_i(x) + \langle g_i, u - x \rangle \Rightarrow \sum_{i=1}^m f_i(u) \geq \sum_{i=1}^m (f_i(x) + \langle g_i, u - x \rangle)$$

$\underset{=f(u)}{\cancel{\sum}}$

$$= \sum_{i=1}^m f_i(x) + \sum_{i=1}^m \langle g_i, u - x \rangle = \sum_{i=1}^m f_i(x) + \langle \sum_{i=1}^m g_i, u - x \rangle = f(x) + \langle \sum_k g_k, u - x \rangle$$

$\underset{\text{מ长时间}}{\cancel{\sum}}$ $\underset{f(x)}{\cancel{\sum}}$

$$\Leftrightarrow f(u) \geq f(x) + \langle \sum_k g_k, u - x \rangle \Rightarrow \sum_k g_k \in \partial f(x)$$

$\cancel{\sum}$

1.2.6

$$\partial(f_1 + f_2) = \partial f_1 + \partial f_2 - \emptyset$$

$\alpha > 0$ סביר $\partial(\alpha \cdot f) = \alpha \cdot \partial f$

$\partial f(x) = \{v \mid f(x) \leq v\}$ סביר, $x \rightarrow$ מושג היפרפלה $f: \mathbb{R}^d \rightarrow \mathbb{R}$ נורם יסוד לוגר

$$\partial\left(\frac{\lambda}{2}\|w\|^2\right) = \lambda\{w, 0\}$$

'סביר' מושג $\|w\|^2$ \rightarrow מינימום

g מינימום - hinge (w, b) סביר $\lambda\{w, 0\}$ 1.2.4 תוגדר ב-1.3.2

$$g_i - \lambda \sum_{j=1}^m g_j \in \partial \sum_{j=1}^m \lambda \{w, b\}$$

$\cancel{\sum}$

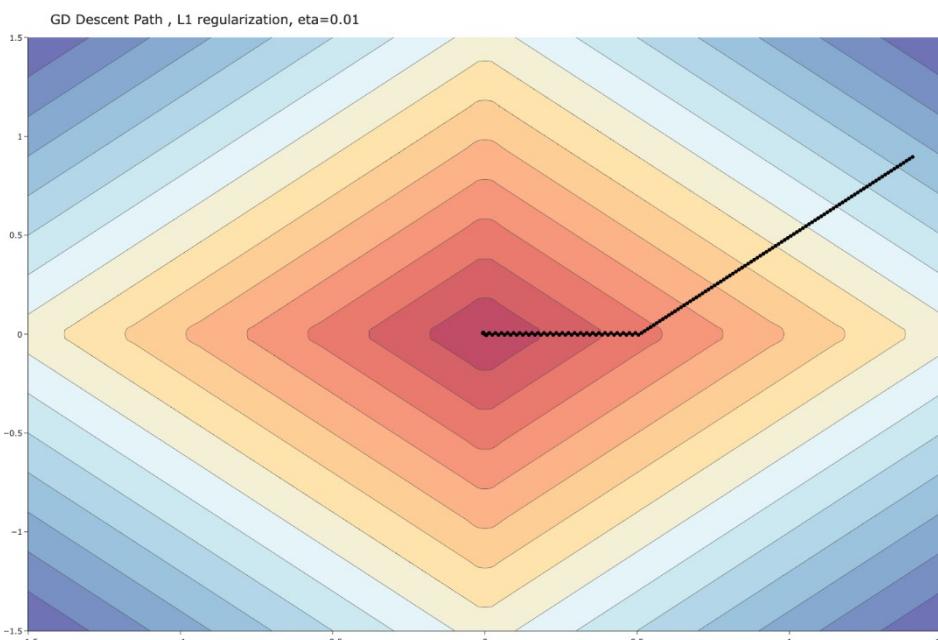
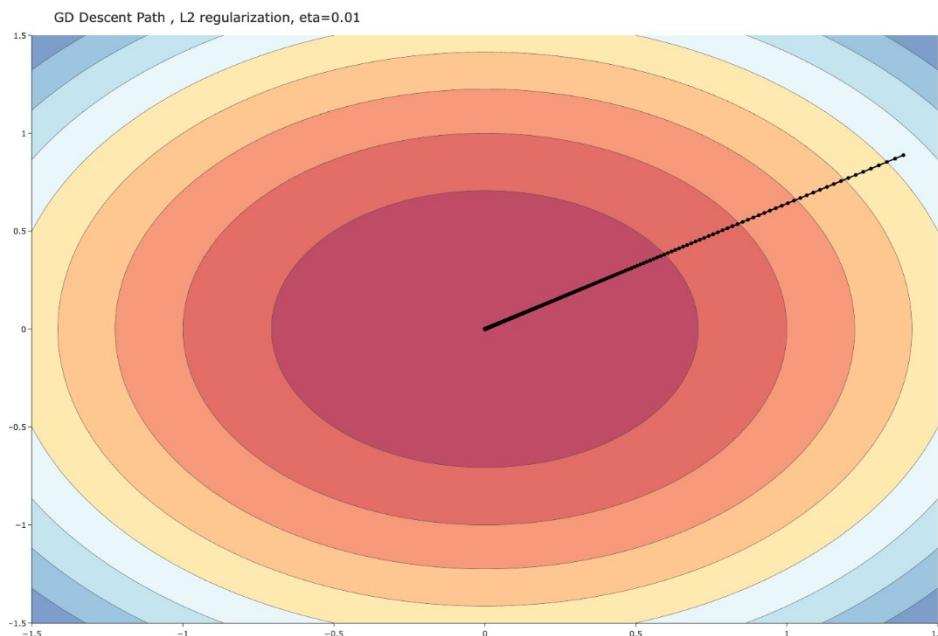
$$\left\{ \lambda(w, b) + \frac{1}{m} \sum_{i=1}^m g_i \mid w \in \mathbb{R}^d, \forall i \in [m], g_i \in \partial \{w, b\} \right\} \in \partial S(w, b)$$

סביר גיאומטריה

Practical part

2.1.1

1

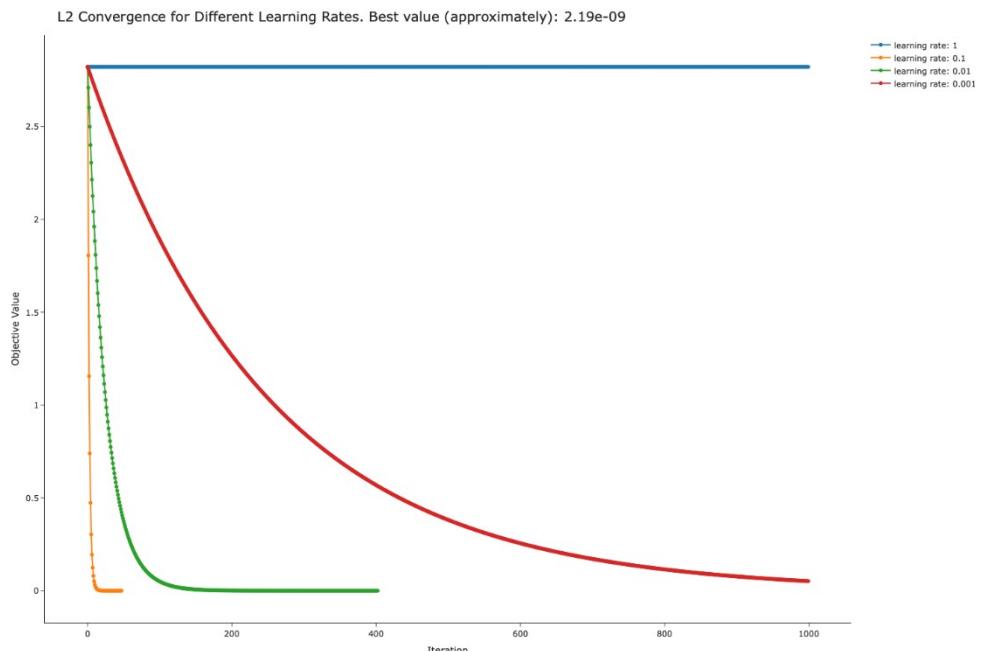
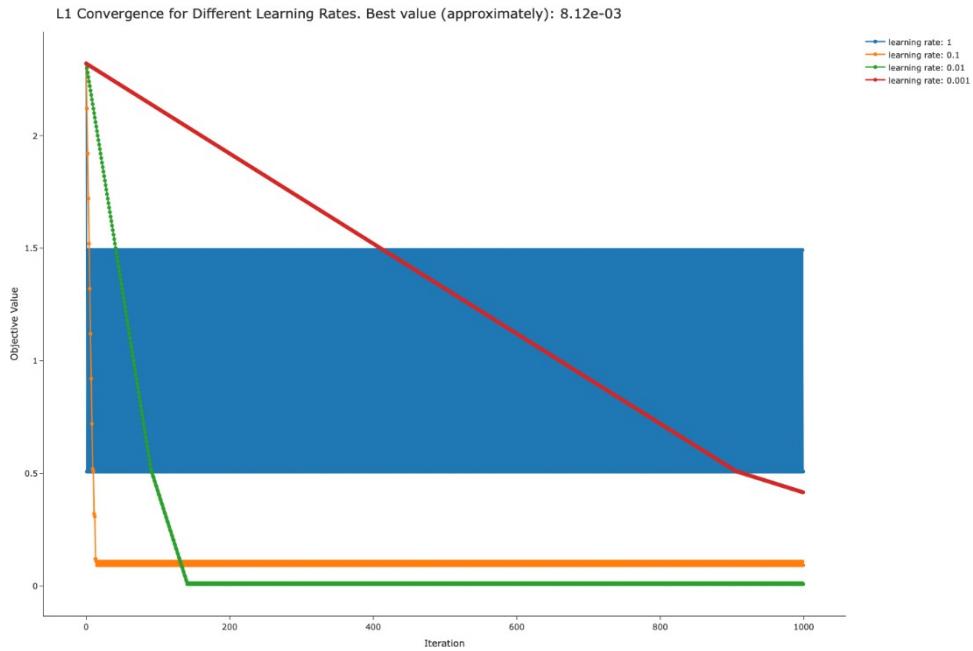


In the L1 objective, the descent path does not follow a straight line toward the minimum, whereas in the L2 objective it does. This difference arises because gradient descent moves in the direction of steepest descent at each step, which is influenced by the geometry of the level sets of the objective. For the L1 objective, level sets are diamond-shaped, so the distance to the minimum varies along the boundary of each level set. As a result, the direction of descent often shifts sharply, causing a more "zig-zag" path. In contrast, the L2 objective has circular level sets centered around the minimum, and the gradient at each point consistently points toward the center, resulting in a smoother and more direct path

2

In the L1 objective, the descent path often changes direction sharply, switching from a diagonal to axis-aligned movement. Near convergence, the algorithm may oscillate between two points due to the non-smoothness of the objective

3



L2 - The convergence curves are smooth across all learning rates, which is expected due to the differentiability and smooth curvature of the L2 objective. When using a learning rate of 1 with the given initialization, gradient descent fails to converge and oscillates between two points. In fact, for any learning rate greater than 1, the algorithm diverges entirely. For learning rates below 1, higher values result in faster convergence

L1 - As with L2, a learning rate of 1 lead to divergence. This happens because the sub gradient at

non-differentiable points in the L1 objective can change direction abruptly, which combined with a large step size causes the algorithm to bounce or overshoot the minimum. The convergence plots for L1 also show “kinks” or sharp bends — a result of the non-smooth geometry of the L1 level sets. Despite this, when convergence does occur, higher learning rates still lead to faster convergence overall.

4

L1 objective best loss $\approx 8.12e-3$

L2 objective best loss $\approx 2.19e-9$

The large gap in final loss values is since the sub gradient of the L1 objective does not shrink as the weights approach the minimum. Combined with a fixed learning rate, this causes gradient descent to stop improving once it gets close enough. In contrast, the gradient of the L2 objective decreases with the weight size, allowing the algorithm to continue making smaller and more precise updates, even with the same fixed step size

5+6+7

