

Theory Questions

Remark: Throughout this exercise, when we write a norm $\|\cdot\|$ we refer to the ℓ_2 -norm.

1. (15 points) Convex functions.

- (a) Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ a convex function, $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$. Show that, $g(\mathbf{x}) = f(A\mathbf{x} + b)$ is convex.

$$\begin{aligned} & \lambda \in [0, 1], \quad w_1, w_2 \in \mathbb{R}^n \quad \text{পৰিস্থিতি} \\ g(\lambda w_1 + (1-\lambda)w_2) & \leq \lambda \cdot g(w_1) + (1-\lambda) \cdot g(w_2) \quad \text{প্ৰমাণ কৰিব দেখা} \\ g(\lambda w_1 + (1-\lambda)w_2) & = f[A(\lambda w_1 + (1-\lambda)w_2) + b] = f(\lambda Aw_1 + (1-\lambda)Aw_2 + b) \\ & = f(\lambda Aw_1 + Aw_2 - \lambda Aw_2 + b) = f(\lambda Aw_1 + \cancel{\lambda b} + Aw_2 + b - \cancel{\lambda Aw_2} - \cancel{\lambda b}) \\ & = f(\lambda(Aw_1 + b) + (1-\lambda)(Aw_2 + b)) \leq \lambda \cdot f(Aw_1 + b) + (1-\lambda) \cdot f(Aw_2 + b) \\ & \quad \text{পৰিস্থিতি} \quad f \text{ is convex} \\ & \quad \Rightarrow \lambda \cdot g(w_1) + (1-\lambda) \cdot g(w_2) \quad \blacksquare \end{aligned}$$

- (b) Consider m convex functions $f_1(\mathbf{x}), \dots, f_m(\mathbf{x})$, where $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$. Now define a new function $g(\mathbf{x}) = \max_i f_i(\mathbf{x})$. Prove that $g(\mathbf{x})$ is a convex function. (Note that from (a) and (b) you can conclude that the hinge loss over linear classifiers is convex.)

$$\begin{aligned} & \lambda \in [0, 1], \quad w_1, w_2 \in \mathbb{R}^n \quad \text{পৰিস্থিতি}, \quad \text{প্ৰমাণ কৰিব দেখা} \\ g(\lambda w_1 + (1-\lambda)w_2) & \leq \lambda \cdot g(w_1) + (1-\lambda) \cdot g(w_2) \quad \text{পৰিস্থিতি} \\ \lambda w_1 + (1-\lambda)w_2 & \text{ দেখা } f_i \rightarrow \text{যদি } f_i \text{ প্ৰমোপণৰ কৰি } f_k \text{ না হৈ } \\ g(\lambda w_1 + (1-\lambda)w_2) & = f_k(\lambda w_1 + (1-\lambda)w_2) \leq \lambda f_k(w_1) + (1-\lambda) \cdot f_k(w_2) = * \\ & \quad \text{পৰিস্থিতি} \quad f_k \text{ is convex} \\ & \quad \forall w \in \mathbb{R}^n, \quad g(w) \geq f_k(w) \quad \text{পৰিস্থিতি} \quad g \text{ নিৰ্দেশ কৰি } g \geq f_k \\ & \quad \Rightarrow * \geq f_k(w_1, w_2) \quad \text{পৰিস্থিতি} \quad g(w_1, w_2) \geq f_k(w_1, w_2) \end{aligned}$$

(c) Let $\ell_{\log} : \mathbb{R} \rightarrow \mathbb{R}$ be the log loss, defined by

$$\ell_{\log}(z) = \log_2(1 + e^{-z})$$

Show that ℓ_{\log} is convex, and conclude that the function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ defined by $f(\mathbf{w}) = \ell_{\log}(y\mathbf{w} \cdot \mathbf{x})$ is convex with respect to \mathbf{w} .

$$\frac{\partial \ell_{\log}(z)}{\partial z} = \frac{-e^{-z}}{(1+e^{-z})\ln(2)}$$

$$\frac{\partial^2 \ell_{\log}(z)}{\partial z^2} = \frac{e^{-z}(1+e^{-z})\ln(2) - (-\ln(2)e^{-z})(-e^{-z})}{[\ln(2)(1+e^{-z})]^2} =$$

$$= \frac{e^{-z}(1+e^{-z})\cdot \ln(2) - e^{-z}\cdot \ln(2)\cdot e^{-z}}{[\ln(2)(1+e^{-z})]^2} = \frac{\ln(2)\cdot e^{-z}\cdot [1+e^{-z}-e^{-z}]}{[\ln(2)(1+e^{-z})]^2}$$

$$= \frac{e^{-z}}{\ln(2)\cdot (1+e^{-z})^2} \Rightarrow \text{for all } z$$

$$\min_{\mathbf{w}} f(\mathbf{w}) = \ell_{\log}(y\mathbf{w} \cdot \mathbf{x})$$

Now we want to show that $f(\mathbf{w})$ is convex. We know that $\ell_{\log}(z)$ is convex. Now we want to show that $\min_{\mathbf{w}} f(\mathbf{w})$ is convex. We know that $\ell_{\log}(z)$ is convex. Now we want to show that $f(\mathbf{w})$ is convex.

4. (15 points) Gradient Descent on Smooth Functions. We say that a continuously differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is β -smooth if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|^2$$

In words, β -smoothness of a function f means that at every point \mathbf{x} , f is upper bounded by a quadratic function which coincides with f at \mathbf{x} .

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a β -smooth and non-negative function (i.e., $f(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$). Consider the (non-stochastic) gradient descent algorithm applied on f with constant step size $\eta > 0$:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$$

Assume that gradient descent is initialized at some point \mathbf{x}_0 . Show that if $\eta < \frac{2}{\beta}$ then

$$\lim_{t \rightarrow \infty} \|\nabla f(\mathbf{x}_t)\| = 0$$

(Hint: Use the smoothness definition with points \mathbf{x}_{t+1} and \mathbf{x}_t to show that $\sum_{t=0}^{\infty} \|\nabla f(\mathbf{x}_t)\|^2 < \infty$ and recall that for a sequence $a_n \geq 0$, $\sum_{n=1}^{\infty} a_n < \infty$ implies $\lim_{n \rightarrow \infty} a_n = 0$. Note that f is not assumed to be convex!)

$$\begin{aligned} & \text{: } \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \text{ 6 wif p"an sl, n } \beta \text{ jno f "ai } n < \frac{2}{\beta} \text{ "i} \\ & f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|^2 \\ & f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^T (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{\beta}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \end{aligned}$$

$$\therefore \text{GD } \Rightarrow \mathbf{x}_{t+1} = \mathbf{x}_t, \quad \text{noj}$$

$$\eta \nabla f(\mathbf{x}_t) = \mathbf{x}_t - \mathbf{x}_{t+1} \iff \mathbf{x}_{t+1} - \mathbf{x}_t = -\eta \nabla f(\mathbf{x}_t) \iff \mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$$

$$\therefore \text{GD } \Rightarrow \mathbf{x}_{t+1} = \mathbf{x}_t \text{ "i} \text{ GD } \Rightarrow \mathbf{x}_{t+1} = \mathbf{x}_t$$

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^T \cdot (-\eta \nabla f(\mathbf{x}_t)) + \frac{\beta}{2} \|\eta \nabla f(\mathbf{x}_t)\|^2$$

\iff

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \eta \|\nabla f(\mathbf{x}_t)\|^2 + \frac{\beta}{2} \cdot \eta^2 \|\nabla f(\mathbf{x}_t)\|^2 = f(\mathbf{x}_t) - \eta \|\nabla f(\mathbf{x}_t)\|^2 \left(1 - \frac{\beta}{2} \eta\right)$$

\iff

$$\eta \|\nabla f(\mathbf{x}_t)\|^2 \left(1 - \frac{\beta}{2} \eta\right) \leq f(\mathbf{x}_t) - f(\mathbf{x}_{t+1})$$

$$\therefore \text{GD } \eta < \frac{2}{\beta} \text{ o "i" 2n "i"}$$

$$\eta < \frac{2}{\beta} \iff \eta \frac{\beta}{2} < 0 \implies 1 - \eta \frac{\beta}{2} > 0$$

$\therefore \text{GD } \eta < \frac{2}{\beta} \text{ o "i" 2n "i"}$

$$\eta \|\nabla f(\mathbf{x}_t)\|^2 \leq f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) \iff \|\nabla f(\mathbf{x}_t)\|^2 \leq \frac{f(\mathbf{x}_t) - f(\mathbf{x}_{t+1})}{\eta - \eta^2 \frac{\beta}{2}}$$

$$\|\nabla f(\mathbf{x}_t)\|^2 \leq \frac{1}{C} (f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}))$$

$$\therefore \text{GD } C = \eta - \eta^2 \frac{\beta}{2} \text{ noj}$$

: min, max "über" über t von Punktfunktionen je nach Projektion

$$\sum_{t=0}^T \|\nabla f(x_t)\|^2 \leq 1 \cdot \sum_{t=0}^T f(x_t) - f(x_{t+1})$$

: $\sum_{t=0}^T \|\nabla f(x_t)\|^2 \leq \frac{1}{C} \cdot (f(x_0) - f(x_{T+1})) \leq \frac{1}{C} \cdot f(x_0) \leq \infty$

$\forall x \in \mathbb{R}^n: f(x) \geq 0$

$$\sum_{t=0}^{\infty} \|\nabla f(x_t)\|^2 \leq \infty \quad \text{für } t \rightarrow \infty \quad \forall t \in \mathbb{N} \quad \text{für } f(x) \geq 0$$

so für alle $t \in \mathbb{N}$ gilt $\|\nabla f(x_t)\|^2 \geq 0$ und

$\lim_{t \rightarrow \infty} \|\nabla f(x_t)\| \rightarrow \infty$: wegen $\sqrt{\sum_{i=1}^d f_i^2(x_t)} \leq \|\nabla f(x_t)\|$ $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\|^2 = \infty$

3. (20 points) GD with projection. In the context of convex optimization, sometimes we would like to limit our solution to a convex set $\mathcal{K} \subseteq \mathbb{R}^d$; that is,

$$\begin{aligned} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in \mathcal{K} \end{aligned}$$

for a convex function f and a convex set \mathcal{K} . In this scenario, each step in the gradient descent algorithm might result in a point outside \mathcal{K} . Therefore, we add an additional projection operator finds the closest point in the set, i.e.:

$$\Pi_{\mathcal{K}}(\mathbf{y}) := \arg \min_{\mathbf{x} \in \mathcal{K}} \|\mathbf{x} - \mathbf{y}\|_2$$

A modified iteration in the *gradient descent with projection* therefore consists of:

$$\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &= \Pi_{\mathcal{K}}(\mathbf{y}_{t+1}) \end{aligned}$$

- (a) Let $\mathbf{y} \in \mathbb{R}^d$ and $\mathbf{x} = \Pi_{\mathcal{K}}(\mathbf{y})$. Prove that for any $\mathbf{z} \in \mathcal{K}$, we have $\|\mathbf{y} - \mathbf{z}\|_2 \geq \|\mathbf{x} - \mathbf{z}\|_2$.
Guidance: use the projection definition and the fact that for any $\lambda \in (0, 1)$, $(1 - \lambda)\mathbf{x} + \lambda\mathbf{z} \in \mathcal{K}$ to show that $\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{z} \rangle \leq \lambda \|\mathbf{z} - \mathbf{x}\|_2^2$ for any $\lambda \in (0, 1)$. Conclude that $\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{z} \rangle \leq 0$. Use that to show the claim in the question.)

וניה גזב \mathcal{K} ולו $x = \Pi_{\mathcal{K}}(y)$, $\lambda \in (0, 1)$, $z \in \mathcal{K}$, $y \in \mathbb{R}^d$

$$(1-\lambda)x + \lambda z \in \mathcal{K} \quad \text{: לפי } x, z \in \mathcal{K}, \text{ אז } x + z \in \mathcal{K}$$

$$x = (1-\lambda)x + \lambda z \quad \text{הו!}$$

בנוסף $y \notin \mathcal{K}$ ומכיוון $x \in \mathcal{K}$ אז $x - y \perp \mathcal{K}$

$$\|x - y\|^2 \leq \|x - z\|^2 \Leftrightarrow \|x - y\| \leq \|x - z\|$$

הו! בזבז לא מוגדר

$$\|x - y\|^2 = \|(1-\lambda)x + \lambda z - y\|^2 = \|x - \lambda x + \lambda z - y\|^2 = \|x - \lambda(x - z) - y\|^2$$

$$= \|x - y - \lambda(x - z)\|^2 = \langle x - y - \lambda(x - z), x - y - \lambda(x - z) \rangle = *$$

$$\therefore \text{for } b = \lambda(x-z), a = x-y \text{ proj of } y$$

$$*= \langle a-b, a-b \rangle = \langle a, a \rangle - 2\langle a, b \rangle + \langle b, b \rangle = \|a\|^2 - 2\langle a, b \rangle + \|b\|^2$$

$\therefore \text{proj of } y \text{ to } \text{span}(x)$

$$*= \|x-y\|^2 - 2\langle x-y, \lambda(x-z) \rangle + \|\lambda(x-z)\|^2$$

$\forall \lambda \in \mathbb{R}$ $a \in \text{span}(x)$ $b = \lambda(x-z)$ $\|a\|^2 - 2\langle a, b \rangle + \|\lambda(x-z)\|^2$

$$-2\langle x-y, \lambda(x-z) \rangle = -2\lambda \langle x-y, x-z \rangle = -2\lambda \langle x-y, x-z \rangle \quad \therefore \text{proj of } y$$

$$*= \|x-y\|^2 - 2\lambda \langle x-y, x-z \rangle + \lambda^2 \|x-z\|^2 \quad \therefore \text{proj of } y$$

$$\|x-y\|^2 \leq \|x-y\|^2 - 2\lambda \langle x-y, x-z \rangle + \lambda^2 \|x-z\|^2 \quad \therefore \text{proj of } y$$

\iff

$$0 \leq -2\lambda \langle x-y, x-z \rangle + \lambda^2 \|x-z\|^2$$

\iff

$$\langle x-y, x-z \rangle \leq \frac{\lambda}{2} \cdot \|x-z\|^2$$

$\forall \lambda \rightarrow 0^+$ $y \in \text{span}(x)$ $\exists z \in K$ $\|y-z\|^2 \geq \|x-z\|^2$

$$\langle x-y, x-z \rangle \leq 0 \quad \therefore \text{proj of } y$$

$$\|y-z\|^2 \geq \|x-z\|^2 \iff \|y-z\| \geq \|x-z\| \quad \therefore \text{proj of } y$$

$$\|y-z\|^2 = \|y-x+x-z\|^2 = \langle a+b, a+b \rangle = \|a\|^2 + 2\langle a, b \rangle + \|b\|^2$$

$$= \|y-x\|^2 + 2\langle y-x, x-z \rangle + \|x-z\|^2 = \|y-x\|^2 - 2\langle x-y, x-z \rangle + \|x-z\|^2$$

$\|y-x\|^2 - 2\langle x-y, x-z \rangle + \|x-z\|^2 \geq 0$ $\therefore \text{proj of } y$

$\therefore \text{proj of } y$

$$\|y-z\|^2 - \|x-z\|^2 \geq 0 \quad \therefore \text{proj of } y$$

$$\|y-z\|^2 - \|x-z\|^2 = \|y-x\|^2 - 2\langle x-y, x-z \rangle + \|x-z\|^2 - \|x-z\|^2 \geq 0$$

(*) $\therefore \text{proj of } y$

(b) Prove that the convergence theorem for GD (as stated in lecture #5) still holds.

$$x^* = w^* \in \arg \min_w f(w), \quad \forall t : x_t := w_t \text{ (no)}$$

We now turn to prove convergence to global minimum for GD applied to a convex and differentiable function. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be such a function, with $w^* \in \arg \min_w f(w)$. Assume that $\|w^*\|_2 \leq B$ and that $\|\nabla f(w)\|_2 \leq G$ for every $w \in \mathbb{R}^d$. We will show that $\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$, the average point over all iterations of GD, converges to the optimal value of $f(\cdot)$.

Let $\epsilon > 0$. For $w_1 = w' = 0$, $\eta_t \equiv \frac{\epsilon}{G^2}$ and $T = \frac{B^2 G^2}{\epsilon^2}$

$$f(\bar{w}) - f(w^*) \leq \epsilon \quad : \text{proof below} \quad \text{jensen inequality}$$

$$f(\bar{w}) = f\left(\frac{1}{T} \sum_{t=1}^T w_t\right) \leq \frac{1}{T} \sum_{t=1}^T f(w_t) \iff f(\bar{w}) - f(w^*) \leq \frac{1}{T} \sum_{t=1}^T f(w_t) - f(w^*)$$

$$= \frac{1}{T} \sum_{t=1}^T (f(w_t) - f(w^*))$$

$$f(\bar{w}) - f(w^*) \leq \frac{1}{T} \sum_{t=1}^T \langle \nabla f(w_t), w_t - w^* \rangle \quad : \text{by rule above}$$

$$f(\bar{w}) - f(w^*) \leq \frac{1}{T} \sum_{t=1}^T \langle u, v \rangle \quad : \text{by } u := w_t - w^*, v := \nabla f(w_t) \text{ no rel w}$$

$$: \|u - \eta v\|^2 \quad \text{if } \eta \in [0, 1]$$

$$\|u - \eta v\|^2 = \langle u - \eta v, u - \eta v \rangle = \|u\|^2 - 2\eta \langle u, v \rangle + \eta^2 \|v\|^2$$

$$: \text{right sign right } \langle u, v \rangle \text{ not } \langle v, u \rangle$$

$$\langle u, v \rangle = \frac{\|u\|^2 + \eta^2 \|v\|^2 - \|u - \eta v\|^2}{2\eta} =$$

$$= \frac{\|w_0 - w^*\|^2 + \eta^2 \|\nabla f(w_0)\|^2 - \|w_t - w^* - \eta \nabla f(w_t)\|^2}{2\eta} \quad : y_{t+1} = w^*$$

$$f(\bar{w}) - f(w^*) \leq \frac{1}{T} \sum_{t=1}^T \langle u, v \rangle = \frac{1}{T} \sum_{t=1}^T \left(\frac{\|w_t - w^*\|^2 - \|y_{t+1} - w^*\|^2}{2\eta} + \frac{\eta}{2} \|\nabla f(w_t)\|^2 \right) \quad : \text{using rule above}$$

$$\|y_{t+1} - x^*\|^2 \geq \frac{1}{2} \|\eta \nabla f(w_t)\|^2 \quad \|\eta \nabla f(w_t)\|^2 \geq \|\nabla f(w_t)\|^2 - \epsilon \quad \text{proof below} \quad \text{proof}$$

$$\therefore (\text{proof point } \|\eta \nabla f(w_t)\|^2 \geq \|\nabla f(w_t)\|^2 - \epsilon \text{ proof}) \quad \|y_{t+1} - x^*\|^2 \geq \frac{1}{2} \|\nabla f(w_t)\|^2 - \frac{\epsilon}{2}$$

• (P) چون پوکن پیش از آن که بگذرد (برای اینجا ممکن است) $\|x_{t+1} - x^*\|^2 \approx$

$$\begin{aligned}
 f(\bar{w}) - f(w^*) &\leq \frac{1}{T} \cdot \frac{1}{2N} \sum_{t=1}^T \|w_t - w^*\|^2 - \|x_{t+1} - w^*\|^2 + \frac{n}{2} \|\nabla f(w_t)\|^2 \\
 &\stackrel{\text{SC-1}}{\leq} \frac{1}{2NT} \left(\|w_1 - w^*\|^2 - \|w_{T+1} - w^*\|^2 \right) + \frac{n}{2} \cdot \frac{1}{T} \sum_{t=1}^T \|\nabla f(w_t)\|^2 \\
 &\stackrel{\text{SC-2}}{\leq} \frac{\|w - w^*\|^2}{2NT} + \frac{n}{2} \cdot \frac{1}{T} \sum_{t=1}^T \|\nabla f(w_t)\|^2 \stackrel{w_t = 0}{\leq} \frac{\|w^*\|^2}{2NT} + \frac{n}{2T} \sum_{t=1}^T \|\nabla f(w_t)\|^2 \\
 &\leq \frac{B^2}{2NT} + \frac{n}{2} \cdot \frac{\epsilon^2}{B^2} = \frac{B^2}{2 \cdot \frac{B^2}{\epsilon^2} + \frac{B^2}{\epsilon^2}} + \frac{\epsilon^2}{2B^2} = \frac{\epsilon^2}{2} + \frac{\epsilon^2}{2} = \epsilon
 \end{aligned}$$

2. (15 points) Hinge loss with linearly separable data. Consider a setup of learning linear classifiers over a dataset $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{-1, +1\}$ for all i . Assume the data is linearly separable, i.e., given a training set there exists $\mathbf{w}^* \in \mathbb{R}^d$ such that $y_i \mathbf{w}^* \cdot \mathbf{x}_i > 0$ for all i (note the strict inequality; assume no bias for simplicity). Recall the definition of the *hinge loss* given in lecture #5:

$$\ell_{\text{hinge}}(r) = \max\{0, 1 - r\}.$$

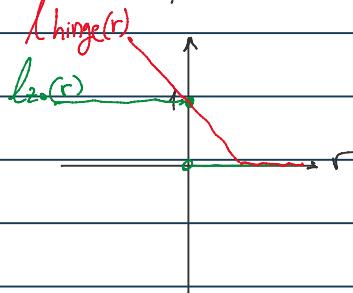
We would like to show that in the linearly separable case, minimizing the hinge loss over the training data will yield a classifier with optimal zero-one loss. Formally, let

$$\mathbf{w}_{\text{hinge}}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ \sum_{i=1}^n \ell_{\text{hinge}}(y_i \mathbf{w} \cdot \mathbf{x}_i) \right\}.$$

Show that $\operatorname{sign}(\mathbf{w}_{\text{hinge}}^* \cdot \mathbf{x}_i) = y_i$ for all i , meaning $\mathbf{w}_{\text{hinge}}^*$ achieves optimal zero-one loss.

(Hint: First show that the hinge loss upper bounds the zero-one loss. Then consider the hinge loss of $c\mathbf{w}^*$ for some constant $c > 0$. What happens to the hinge loss as $c \rightarrow \infty$?)

$\ell_{\text{hinge}}(r) \geq \ell_{\text{zo}}(r) : r \leq 1 \Rightarrow \text{zo-loss} \leq \text{hinge-loss}$



$\forall i: y_i \langle w^*, x_i \rangle \geq 0$ if w^* is linearly separable (i.e. $y_i \geq 0$)

↳ even if L_2 loss hinge w^* has large slope

: (1) if hinge $c > 0$ then $\|w^*\|_2$ is large

$$l_{\text{hinge}}(y_i \langle c w^*, x_i \rangle) = \max\{0, 1 - y_i \langle c w^*, x_i \rangle\} = \max\{0, 1 - c \cdot y_i \langle w^*, x_i \rangle\}$$

, if $c < y_i \langle w^*, x_i \rangle$ a small w^* has large $\|w^*\|_2$

: if c large w^* has large $\|w^*\|_2$

$$\lim_{c \rightarrow \infty} l_{\text{hinge}}(y_i \langle c w^*, x_i \rangle) = \lim_{c \rightarrow \infty} \max\{0, 1 - c \cdot y_i \langle w^*, x_i \rangle\} = \max\{0, -\infty\} = 0$$

then it is good for large hinge-loss if c is large and w^* is small

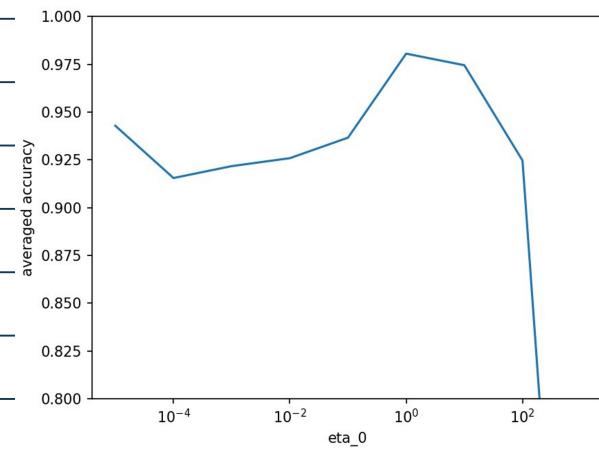
hinge loss is zero if w^* is large and L_2 loss hinge is zero

: if w^* is large hinge loss is zero. If L_2 loss hinge is zero then w^* is large

* if w^* is large $\Rightarrow \text{sign}(w^* \langle w^*, x_i \rangle) = y_i$

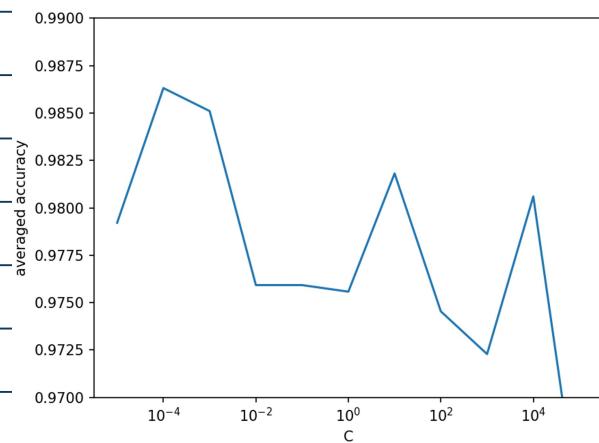
accuracy

SGD for hinge loss - 1 noise



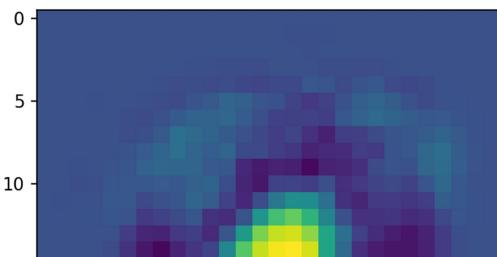
(a)

accuracy vs C no noise

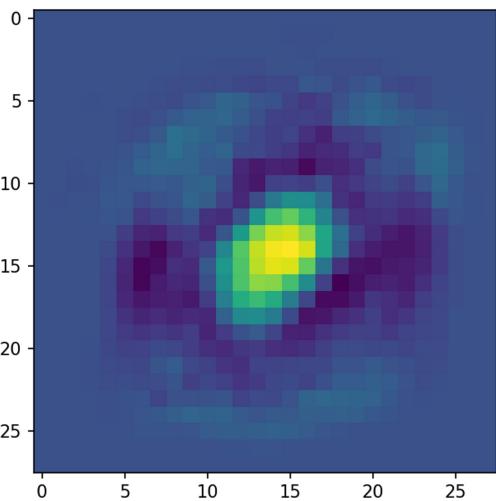


(b)

accuracy vs C no noise



(c)



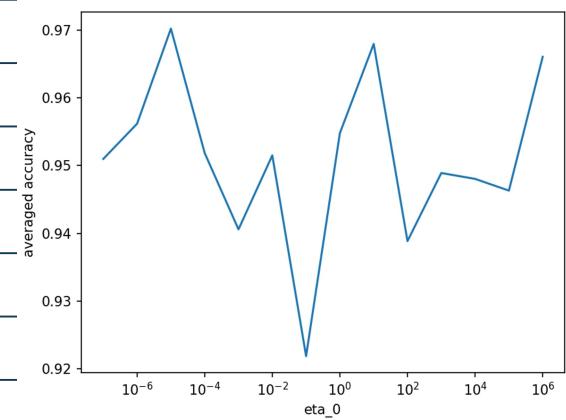
(c)

• 7/8 points for correctly classifying the digit '4' with accuracy 0.9928352098259979

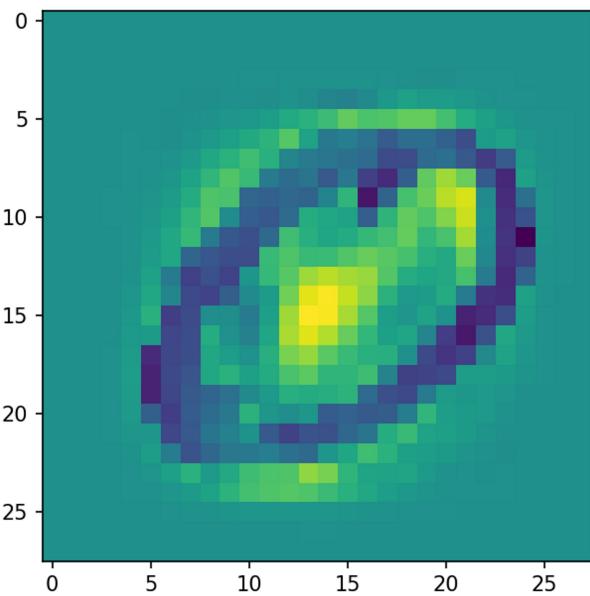
the accuracy of the best classifier on the test set is: 0.9928352098259979

(d)

SGD for log-loss - 2 nice

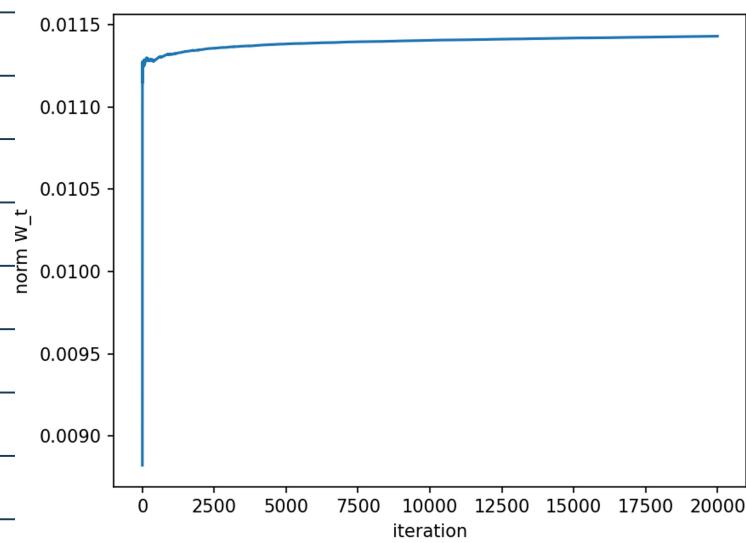


(e)



(b)

the accuracy of the best classifier on the test set is: 0.9646878198567042



(c)

המבחן מוצג כטבלה הבאה. בזאת, ניתן לרשום את המשוואת הדרישה. מטרת הימינית היא $\|w\|_2^2$. מטרת המינימיזציה היא $\sum_{i=1}^{1000} \max(0, 1 - y_i(w^\top x_i + b))$.