

## Theory Questions

1. (15 points) **SVM with multiple classes.** One limitation of the standard SVM is that it can only handle binary classification. Here is one extension to handle multiple classes. Let  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  and now let  $y_1, \dots, y_n \in [K]$ , where  $[K] = \{1, 2, \dots, K\}$ . We will find a separate classifier  $\mathbf{w}_j$  for each one of the classes  $j \in [K]$ , and we will focus on the case of no bias ( $b = 0$ ). Define the following loss function (known as the *multiclass hinge-loss*):

$$\ell(\mathbf{w}_1, \dots, \mathbf{w}_K, \mathbf{x}_i, y_i) = \max_{j \in [K]} (\mathbf{w}_j \cdot \mathbf{x}_i - \mathbf{w}_{y_i} \cdot \mathbf{x}_i + \mathbb{1}(j \neq y_i)),$$

where  $\mathbb{1}(\cdot)$  denotes the indicator function. Define the following multiclass SVM problem:

$$f(\mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}_1, \dots, \mathbf{w}_K, \mathbf{x}_i, y_i)$$

After learning all the  $\mathbf{w}_j$ ,  $j \in [K]$ , classification of a new point  $\mathbf{x}$  is done by  $\arg \max_{j \in [K]} \mathbf{w}_j \cdot \mathbf{x}$ . The rationale of the loss function is that we want the "score" of the true label,  $\mathbf{w}_{y_i} \cdot \mathbf{x}_i$ , to be larger by at least 1 than the "score" of each other label,  $\mathbf{w}_j \cdot \mathbf{x}_i$ . Therefore, we pay a loss if  $\mathbf{w}_{y_i} \cdot \mathbf{x}_i - \mathbf{w}_j \cdot \mathbf{x}_i \leq 1$ , for  $j \neq y_i$ .

Consider the case where the data is linearly separable. Namely, there exists  $\mathbf{w}_1^*, \dots, \mathbf{w}_K^*$  such that  $y_i = \arg \max_j \mathbf{w}_j^* \cdot \mathbf{x}_i$  for all  $i$ . Show that any minimizer of  $f(\mathbf{w}_1, \dots, \mathbf{w}_K)$  will have zero classification error.

פונקציית  $f(\mathbf{w}_1, \dots, \mathbf{w}_K)$  מגדילה נזק לאנרגיה כיווןclassification error=0 ופונקציית האנרגיה מינימלית.

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_K} f(\mathbf{w}_1, \dots, \mathbf{w}_K) = \min \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}_1, \dots, \mathbf{w}_K, \mathbf{x}_i, y_i)$$

$$= \min_{\mathbf{w}_1, \dots, \mathbf{w}_K} \frac{1}{n} \sum_{i=1}^n \max_{j \in [K]} (\mathbf{w}_j \cdot \mathbf{x}_i - \mathbf{w}_{y_i} \cdot \mathbf{x}_i + \mathbb{1}(j \neq y_i))$$

• מינימיזציה של אנרגיה כיוון  $j=y_i$  וזרה מינימלית

$$\mathbf{w}_j \cdot \mathbf{x}_i - \mathbf{w}_j \cdot \mathbf{x}_i + 0 = 0$$

כדי ש  $y_j \neq y_i$  ובדגש נזכיר. קטע ה-

$$w_j \cdot x_i - w_{y_i} \cdot x_i + 1$$

לפניך רצויים

$$\min_{w_1, \dots, w_K} \frac{1}{n} \sum_{i=1}^n \left( \max_{\substack{j \in [K] \\ j \neq y_i}} \{w_j \cdot x_i - w_{y_i} \cdot x_i + 1, 0\} \right)$$

כונך פולג כ'

$$\max_{\substack{j \in [K] \\ j \neq y_i}} \{w_j \cdot x_i - w_{y_i} \cdot x_i + 1, 0\}$$

ולפניך נשים  $w_1, \dots, w_K$  כפופה!

$$\max_{\substack{j \in [K] \\ j \neq y_i}} \{w_j \cdot x_i - w_{y_i} \cdot x_i + 1, 0\} \leq 0$$

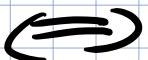
ונשים  $\max$  זה ותנו "

$$\max_{\substack{j \in [K] \\ j \neq y_i}} (w_j \cdot x_i - w_{y_i} \cdot x_i + 1) \leq 0$$

בנוסף ליותר, רצינו ש



$$\forall j \in [K], j \neq y_i \quad w_j \cdot x_i - w_{y_i} \cdot x_i + 1 \leq 0$$



$$\forall j \in [K], j \neq y_i \quad w_j \cdot x_i - w_{y_i} \cdot x_i \leq -1$$



$$\forall j \in [K], j \neq y_i \quad w_{y_i} \cdot x_i - w_j \cdot x_i \geq 1$$

כונך ה-  $w_1, \dots, w_K$  מושג יי'ר

פוא'

למיינר פונקציית האנרגיה מוגדרת כ

$$y_i^* = \arg \max_y w_y^* \cdot x_i \text{ ו } w_1^*, \dots, w_k^*$$

- $w_{y_i^*} \cdot x_i > w_j \cdot x_i \text{ ו } y_i^* \neq j \in [k]$
- $w_{y_i^*} \cdot x_i - w_j \cdot x_i > 0$

רעיון:

$$\text{, } 1 \leq i \leq n \text{ מתי } w_k^* = \frac{1}{\min_{i,j \neq y_i^*} w_{y_i^*} \cdot x_i - w_j \cdot x_i} \cdot w_k^*$$

: מ"מ  $j \neq y_i^*$

$$w_{y_i^*} \cdot x_i - w_j \cdot x_i = \frac{w_{y_i^*} \cdot x_i}{\min_{i,j \neq y_i^*} w_{y_i^*} \cdot x_i - w_j \cdot x_i}$$

$$- \frac{w_j \cdot x_i}{\min_{i,j \neq y_i^*} w_{y_i^*} \cdot x_i - w_j \cdot x_i} = \frac{w_{y_i^*} \cdot x_i - w_j \cdot x_i}{\min_{i,j \neq y_i^*} w_{y_i^*} \cdot x_i - w_j \cdot x_i} \geq 1 \quad \downarrow (*)$$

(\*) הוכחה: כرار הנעה הינה  
כך ש  $y_i^*$  מושג בדיעבד מה  $y_i$  מושג בדיעבד  
ולכן  $w_{y_i^*} \cdot x_i - w_j \cdot x_i > 0$

כרא ריגענו נניח לא  $w_1, \dots, w_k$  כדלהלן:

$$\forall j \in [k], j \neq y_i \quad w_{y_i} \cdot x_i - w_j \cdot x_i \geq 1$$

: ופ"ז

$$\max_{[k]} \{ w_j \cdot x_i - w_{y_i} \cdot x_i + 1, 0 \} \leq 0$$

כך

$$\max_{[K]} \{w_j \cdot x_i - w_{j+1} \cdot x_{i+1}, 0\} \geq 0$$

כל קאכין:

$$\max_{[K]} \{w_j \cdot x_i - w_{j+1} \cdot x_{i+1}, 0\} = 0$$

$$\min_{w_1, \dots, w_K} f(w_1, \dots, w_K) = 0 \quad \text{כלומר}$$

רנו  $w_1, \dots, w_K$  כך ש  $f(w_1, \dots, w_K) \geq 0$  וענוק  $w_1, \dots, w_K$  כך ש  $f(w_1, \dots, w_K) = 0$

$f(w_1, \dots, w_K) = 0$  (הו. נרמול ומיינדר)

$$\min_{w_1, \dots, w_K} f(w_1, \dots, w_K) = 0 \quad \text{כלומר}$$

. classification error = 0

: ש  $f(w_1, \dots, w_K) = 0$  ומי

$$f(w_1, \dots, w_K) = \frac{1}{n} \sum_{i=1}^n \max_{[K]} \{w_j \cdot x_i - w_{j+1} \cdot x_{i+1}, 0\} = 0$$

ומי  $\forall i$   $w_j \cdot x_i - w_{j+1} \cdot x_{i+1} \geq 0$

לעומת רגולציה כפולה נתקע

:  $1 \leq i \leq n$   $\exists j$   $w_j \cdot x_i - w_{j+1} \cdot x_{i+1} > 0$

$$\max_{[K]} \{w_j \cdot x_i - w_{j+1} \cdot x_{i+1}, 0\} = 0$$

כלומר:

$$w_j \cdot x_i - w_{j+1} \cdot x_{i+1} \leq 0 \Rightarrow w_j \cdot x_i - w_{j+1} \cdot x_i \leq -1$$

$$\Rightarrow w_{j+1} \cdot x_i - w_j \cdot x_i \geq 1 \Rightarrow w_{j+1} \cdot x_i \geq w_j \cdot x_i + 1$$

$$\Rightarrow w_{j+1} \cdot x_i > w_j \cdot x_i \quad \forall 1 \leq i \leq n$$

$$\underset{j \in [k]}{\operatorname{argmax}} w_j \cdot x_i = y_i \quad : \text{פוד}$$

כדי שפונקציית הטעינה תהיה מוגדרת כפונקציית קבוצה  
classification error=0

2. (10 points) Expressivity of ReLU networks. Consider the ReLU activation function:

$$h(x) = \max\{x, 0\}$$

Show that the maximum function  $f(x_1, x_2) = \max\{x_1, x_2\}$  can be implemented using a neural network with one hidden layer and ReLU activations. You can assume that there is no activation after the last layer.

(Hint: It is possible to implement the function using a hidden layer with 4 neurons. You may find the following identity useful:  $\max\{x_1, x_2\} = \frac{x_1+x_2}{2} + \frac{|x_1-x_2|}{2}$ .)

$$\therefore \text{SK } x \geq 0 \text{ PK}$$

$$x = \max\{x, 0\} = h(x), \quad 0 = \max\{-x, 0\} = h(-x)$$

$$x - x - 0 = \max\{x, 0\} - \max\{-x, 0\}$$

$$(x \geq 0 \text{ נס סימול}): \text{לפניך}$$

$$x \geq -x \Rightarrow |x| = x = \max\{x, 0\} + \max\{-x, 0\} = x + 0 = x$$

$$\therefore \text{SK } x < 0 \text{ PK}$$

$$-x = \max\{-x, 0\} = h(-x), \quad 0 = \max\{x, 0\} = h(x)$$

$$0 = x - x = x + \max\{-x, 0\} = \max\{x, 0\}$$

$$\Rightarrow x = \max\{x, 0\} - \max\{-x, 0\}$$

$$(x < 0 \text{ נס סימול}): \text{לפניך}$$

$$x < -x \Rightarrow |x| = -x = \max\{x, 0\} + \max\{-x, 0\} = 0 - x = -x$$

$$\therefore x \text{ סימול נון}$$

$$x = \max\{x, 0\} - \max\{-x, 0\} = h(x) - h(-x)$$

$$|x| = \max\{x, -x\} = \max\{x, 0\} + \max\{-x, 0\} = h(x) + h(-x)$$

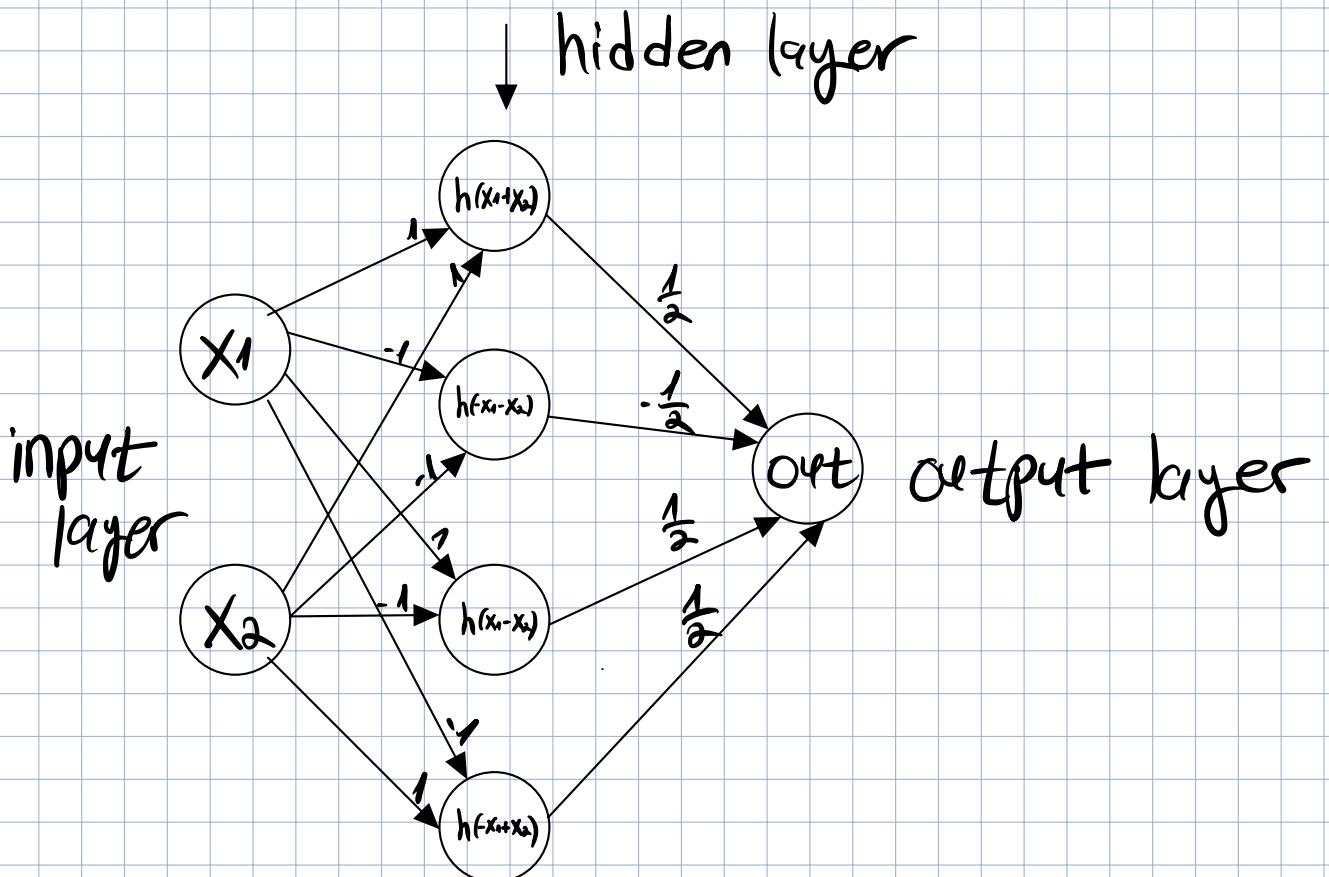
NCCN:

$$\max\{x_1, x_2\} = \frac{x_1 + x_2}{2} + \frac{|x_1 - x_2|}{2}$$

$$\frac{x_1 + x_2}{2} + \frac{|x_1 - x_2|}{2} = h(x_1 + x_2) - h(-x_1 - x_2) + \frac{h(x_1 - x_2) + h(-x_1 + x_2)}{2}$$

$$= \frac{1}{2} [h(x_1 + x_2) - h(-x_1 - x_2) + h(x_1 - x_2) + h(-x_1 + x_2)]$$

לפנינו מושג עליון גראDED:



3. (15 points) Soft SVM with  $\ell^2$  penalty. Consider the following problem:

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n} & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 \\ \text{s.t. } & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n \end{aligned}$$

- (a) Show that a constraint of the form  $\xi_i \geq 0$  will not change the problem. Meaning, show that these non-negativity constraints can be removed. That is, show that the optimal value of the objective will be the same whether or not these constraints are present.
- (b) What is the Lagrangian of this problem?
- (c) Minimize the Lagrangian with respect to  $\mathbf{w}, b, \xi$  by setting the derivative with respect to these variables to 0.
- (d) What is the dual problem?

הוכיח שproblem קבוי מינימום הינו שווה problem שמיון מינימום הינו שווה.  $V_1 = V_2$ .  $V_2 = \text{minimum}_{\mathbf{w}, b, \xi} \text{subject to } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n$

הוכיח שמיון מינימום הינו שווה problem שמיון מינימום הינו שווה.  $V_1 = V_2$ .  $V_1 = \text{minimum}_{\mathbf{w}, b, \xi} \text{subject to } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n$

הוכיח שמיון מינימום הינו שווה problem שמיון מינימום הינו שווה.  $V_1 = V_2$ .  $V_1 = \text{minimum}_{\mathbf{w}, b, \xi} \text{subject to } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n$

הוכיח שמיון מינימום הינו שווה problem שמיון מינימום הינו שווה.  $V_1 = V_2$ .  $V_1 = \text{minimum}_{\mathbf{w}, b, \xi} \text{subject to } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n$

הוכיח שמיון מינימום הינו שווה problem שמיון מינימום הינו שווה.  $V_1 = V_2$ .  $V_1 = \text{minimum}_{\mathbf{w}, b, \xi} \text{subject to } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, n$

$$V_1 = \min_{\substack{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n}} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{2} \sum_{i=1}^n \xi_i^2 = \min_{\substack{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n}} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{2} \cdot \xi_j^2 + \frac{C}{2} \sum_{j \neq i} \xi_i^2$$

(\*)

הוכיח שמיון מינימום הינו שווה problem שמיון מינימום הינו שווה.  $\xi_j = 0$  אם  $y_j (\mathbf{w}^T \mathbf{x}_j + b) \geq 1$  ו- $\xi_j = 1$  אחרת.  $y_j (\mathbf{w}^T \mathbf{x}_j + b) \geq 1 - \xi_j$  אם  $y_j (\mathbf{w}^T \mathbf{x}_j + b) < 1$ .

כדו' נאמר כי אם  $y_j = 1$  אז  $w^\top x_j + b \geq 1 - \epsilon_j$   
ובנוסף לכך  $y_j = 0$  אז  $w^\top x_j + b \leq -1 + \epsilon_j$ .

וונא לשים  $\epsilon_j = 0$  ו/or  $w^\top x_j + b = 1$  הינה גורם אחד  
בהתאם להינתן. וק"מ  $0 = y_j w^\top x_j + b$  מוגדרת גזירה ב  $x_j$   
 $v_1 = v_2$  מוגדרת  $v_2 \leq v_1 \leq v_2$  כוונתית.  $v_1 \geq v_2$   
כלומר הוכחה ב. כי גזירה אחורית טענית.

(b)

$$\min_{w \in \mathbb{R}^n, b \in \mathbb{R}, \varepsilon \in \mathbb{R}^n} \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^n \varepsilon_i^2$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i \quad \forall i = 1, \dots, n$$

האם ימ' אם אומרים כי הרכ麿ת צור כל'ו ניכר בפ'

$$\min_{w \in \mathbb{R}^n, b \in \mathbb{R}, \varepsilon \in \mathbb{R}^n} \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^n \varepsilon_i^2$$

$$\text{s.t. } -y_i(w^T x_i + b) + 1 - \varepsilon_i \leq 0 \quad \forall i = 1, \dots, n$$

: מינ' פ'

$$\mathcal{L}(w, \varepsilon, \alpha, b) =$$

$$= \frac{1}{2} \|w\|^2 + \frac{C}{2} \cdot \sum_{i=1}^n \varepsilon_i^2 + \sum_{i=1}^n \alpha_i (-y_i(w^T x_i + b) + 1 - \varepsilon_i)$$

$$= \frac{1}{2} \|w\|^2 + \frac{C}{2} \cdot \|\varepsilon\|^2 + \sum_{i=1}^n -\alpha_i \varepsilon_i + \alpha_i (-y_i(w^T x_i + b) + 1)$$

$$= \frac{1}{2} \|w\|^2 + \frac{C}{2} \cdot \|\varepsilon\|^2 - \alpha \cdot \varepsilon + \sum_{i=1}^n \alpha_i (-y_i(w^T x_i + b) + 1)$$

(C) מינימיזציה של היגוי ופונקציית האפסה (0-ג'רליאן)

$$\nabla_w \lambda = w - \sum_{i=1}^n \alpha_i y_i x_i = 0 \Leftrightarrow w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\nabla_\varepsilon \lambda = \frac{\partial C}{\partial \varepsilon} \varepsilon - \lambda = C\varepsilon - \lambda \Leftrightarrow C\varepsilon = \lambda \Leftrightarrow \varepsilon = \frac{\lambda}{C}$$

$$\nabla_b \lambda = - \sum_{i=1}^n \alpha_i y_i = 0$$

כדי שתהיה כוונת:

$$\frac{1}{2} \left( \sum_{i=1}^n \alpha_i y_i x_i \right)^2 + \frac{C}{2} \left\| \frac{\lambda}{C} \right\|^2 - \lambda \cdot \frac{\lambda}{C} + \sum_{i=1}^n \alpha_i [-y_i (\sum_{j=1}^n \alpha_j y_j x_j)^T x_i + b] + 1$$

$$= \frac{1}{2} \left( \sum_{i=1}^n \alpha_i y_i x_i \right)^2 + \frac{1}{2C} \|\lambda\|^2 - \frac{1}{2} \frac{\lambda^2}{C} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (\sum_{j=1}^n \alpha_j y_j x_j)^T x_i - \frac{1}{2} \sum_{i=1}^n \alpha_i y_i$$

: מטריה  $\sum_{i=1}^n \alpha_i y_i = 0$  נסכימה

$$\frac{1}{2} \left( \sum_{i=1}^n \alpha_i y_i x_i \right)^2 - \frac{1}{2C} \|\lambda\|^2 + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (\sum_{j=1}^n \alpha_j y_j x_j)$$

$$= -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \cdot x_i \cdot x_j + \sum_{i=1}^n \alpha_i - \frac{1}{2C} \|\lambda\|^2$$

D) מינימיז:

$$\max_{\alpha} \min_{w, \varepsilon, b} L(w, \varepsilon, \alpha, b) = \max_{\alpha} g(\alpha) = \min_{\alpha} -g(\alpha)$$

לינארית (גראף):

$$\begin{aligned} \min_{\alpha} & \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \cdot x_i \cdot x_j - \sum_{i=1}^n \alpha_i + \frac{1}{2C} \cdot \|\alpha\|^2 \\ & - \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

4. (15 points) Gradient of cross-entropy loss over softmax. Let  $\mathbf{y} \in \{0, 1\}^d$  be a one-hot vector, i.e.  $\mathbf{y}$  has a single entry which is 1 and the rest are 0. Consider the following loss function  $\ell_{\mathbf{y}} : \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$\ell_{\mathbf{y}}(\mathbf{w}) = -\mathbf{y} \cdot \log(\text{softmax}(\mathbf{w})),$$

where  $\text{softmax}(\mathbf{w}) = \frac{e^{\mathbf{w}}}{\sum_{j=1}^d e^{w_j}}$ , is the softmax function (see slide 7 of lecture #7). In the above notation, the exponent and logarithm function operate elementwise when given a vector as an input. The function  $\ell_{\mathbf{y}}$  is known as the *cross-entropy loss*, and you will encounter it in the programming assignment.

Prove that the gradient of  $\ell_{\mathbf{y}}$  with respect to  $\mathbf{w}$  is given by:

$$\nabla \ell_{\mathbf{y}}(\mathbf{w}) = \text{softmax}(\mathbf{w}) - \mathbf{y}$$

(רוחיג כהמקלה:

$$\forall 1 \leq i \leq d \quad \frac{\partial \ell_{\mathbf{y}}(\mathbf{w})}{\partial w_i} = -y_i \cdot \frac{1}{\sum_{j=1}^d e^{w_j} \cdot \ln e} \cdot e^{w_i} \cdot \frac{\sum_{j=1}^d e^{w_j} - e^{w_i}}{\left(\sum_{j=1}^d e^{w_j}\right)^2}$$

$$= -y_i \cdot \frac{\sum_{j=1}^d e^{w_j}}{e^{w_i}} \cdot \frac{e^{w_i} \left( \sum_{j=1}^d e^{w_j} - e^{w_i} \right)}{\left( \sum_{j=1}^d e^{w_j} \right)^2}$$

$$= -y_i \cdot \frac{\sum_{j=1}^d e^{w_j} - e^{w_i}}{\sum_{j=1}^d e^{w_j}} = -y_i \left( 1 - \frac{e^{w_i}}{\sum_{j=1}^d e^{w_j}} \right)$$

$$= -y_i + y_i \cdot \frac{e^{w_i}}{\sum_{j=1}^d e^{w_j}}$$

נימוק מהר ורשמי יותר יופיע מאוחר יותר

$$y_i \frac{e^{w_i}}{\sum_{j=1}^c e^{w_j}} = \begin{cases} \frac{e^{w_i}}{\sum_{j=1}^c e^{w_j}} & i=c \\ 0 & \text{if } i \neq c \end{cases}$$

: δρ j ع

$$\nabla \ell_y(w) = \left( \frac{\partial \ell_y}{\partial w_1}, \dots, \frac{\partial \ell_y}{\partial w_d} \right)$$

$$= \left( \frac{e^{w_1}}{\sum_{j=1}^d e^{w_j}} - y_i, \dots, \frac{e^{w_d}}{\sum_{j=1}^d e^{w_j}} - y_i \right) = \text{softmax}(w) - y$$

لرگز

לען כוכבים:

.91% - 1% גמלו

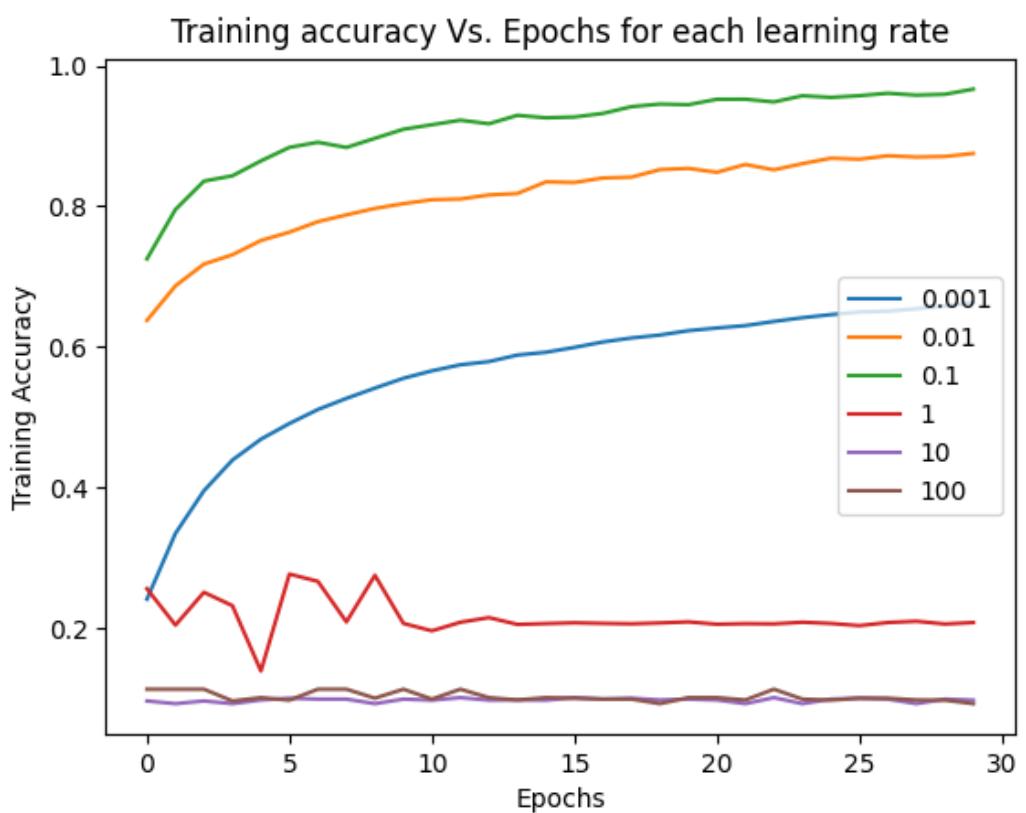
: 1b note

- (b) (10 points) Train a one-hidden layer neural network as in the example given above (e.g., training set of size 10000, one hidden layer of size 40). For each learning rate in  $\{0.001, 0.01, 0.1, 1, 10, 100\}$ , plot the *training* accuracy, *training loss* ( $\ell(\mathcal{W})$ ) and *test* accuracy across epochs (3 plots: each contains the curves for all learning rates). For the test accuracy you can use the *one\_label\_accuracy* function, for the training accuracy use the *one\_hot\_accuracy* function and for the training loss you can use the *loss* function. All functions are in the *Network* class.

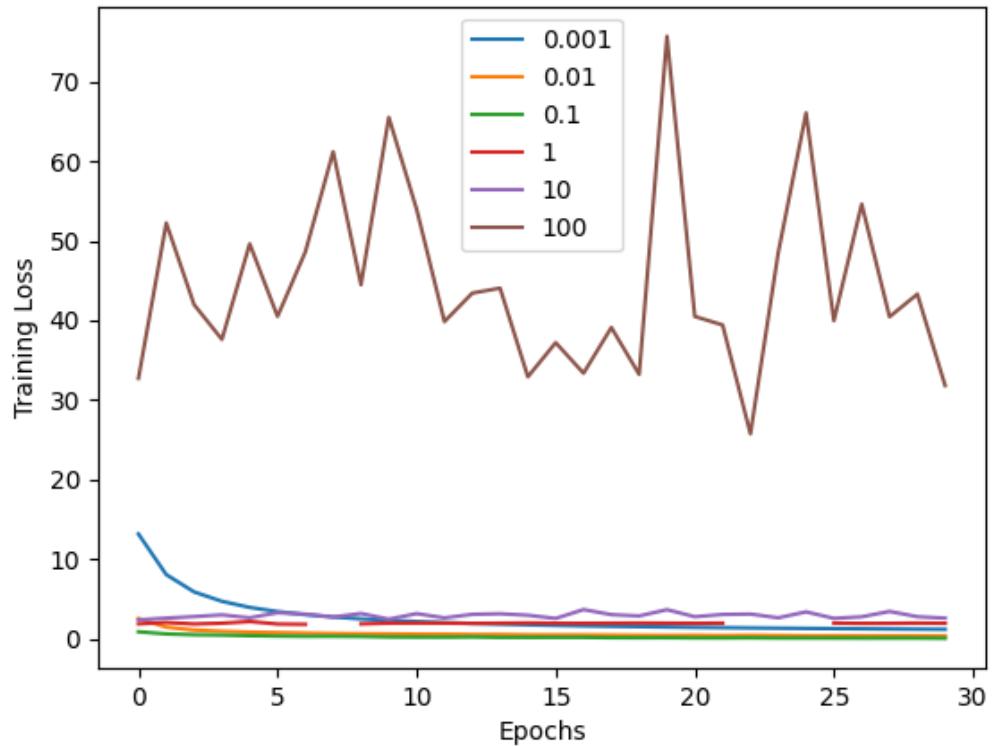
The test accuracy with leaning rate 0.1 in the final epoch should be above 80%.

What happens when the learning rate is too small or too large? Explain the phenomenon.

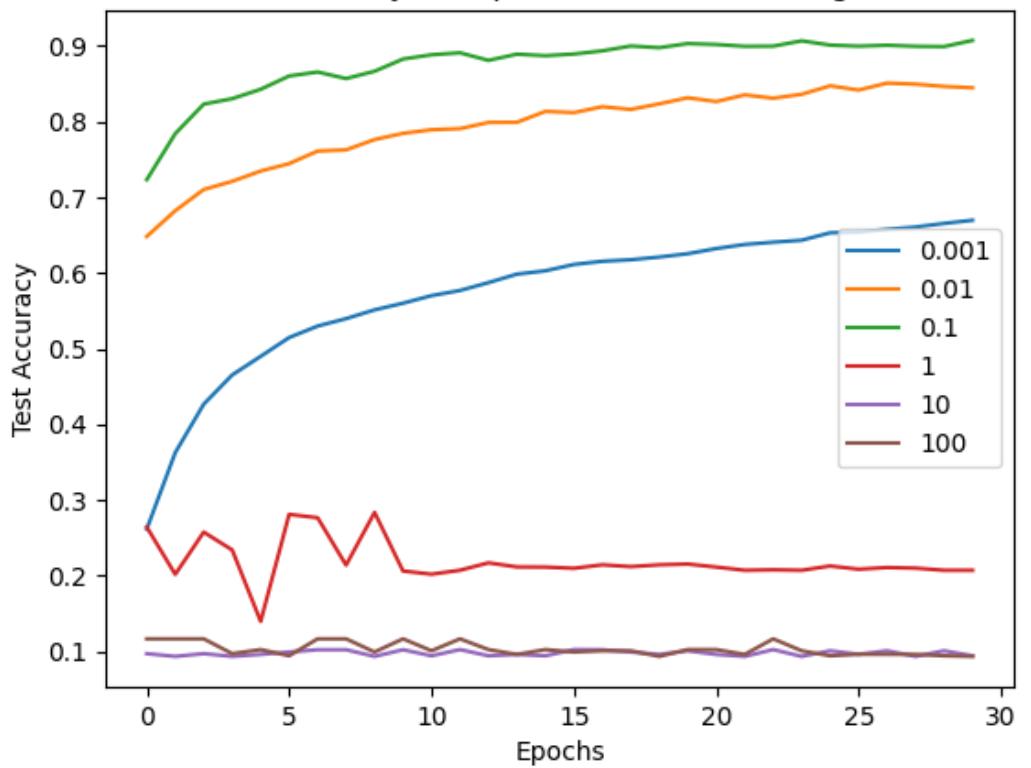
: פותח



Training Loss Vs. Epochs for each learning rate



Test Accuracy Vs. Epochs for each learning rate



כין מודל עם training accuracy ו- test accuracy

הוכין סותי עם סט נתונים (הסטטיסטיקה הוויזואלית).

וכו ו' learning rate=0.1 ו-SGD ייו' גודן.

נב. עליך "זרעו יותר צלחות" נאך זה יצליח?

(אפקט כתם וכחון). ריתן מודולר פונקציית פג'ר  $f(x) = 0.001x^2$ .

ש. ו' בוקס קאפס פלטיאז'ה ווען פלטאליג'ה צו' פונקציית הינה

וגודל learning rate'ה זיך.

וכו ו' learning rate'ה שווה נב. סול'ה ה- $\theta_0$  SGD צו' גודל.

אלאה מכך לדוגמא את פונקטר הוייריאנט, פלאן נורמל נב. מוקדם.

ונדריך (ונילך): (אפקט פלטיאז'ה וווען).

ריתן מודולר פונקציית פונקציית הינה

דינמיות ו' וווען פלטאליג'ה (כוא כה'ן).

(c) (5 points) Now train the network on the whole training set and test on the whole test set:

```
>>> training_data, test_data = data.load(train_size=50000,test_size=10000)
>>> net = network.Network([784, 40, 10])
>>> net.SGD(training_data, epochs=30, mini_batch_size=10, learning_rate=0.1,
test_data=test_data)
```

Do **not** calculate the training accuracy and training loss as in the previous section (this is time consuming). What is the test accuracy in the final epoch (should be above 90%)?

גראן 94.23%

Epoch 29 test accuracy: 0.9423

וילגס

... אנו מודים לך על הצלחה בEpoch 29!

(d) (10 points bonus) Explore different structures and parameters and try to improve the test accuracy. Use the whole test set. In order to get full credit you should get accuracy

of more than 96%. Any improvement over the previous clauses will give you 5 points.

The maximum score for this homework set remains 100.

```
# Q1d - bonus
training_data, test_data = backprop_data.load(train_size=50000, test_size=10000)
net = backprop_network.Network([784, 800, 10])
net.SGD(training_data, 30, 10, 0.1, test_data, 1)
```

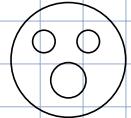
: 91%

Epoch 27 test accuracy: 0.9617

Epoch 28 test accuracy: 0.9617

Epoch 29 test accuracy: 0.9617

הנחיין שפתקה!



! 96.17% גראן!

2.10

- (a) (5 points) Train the network **on 10% of the training set** (otherwise the training process would take a long time) once with the default parameters given in the code, and then again using the Adam optimizer instead of SGD (see the documentation in `torch.optim.Adam` for more information), with a learning rate of 0.00005. How is the training process affected by this change? Do you achieve better training loss? What about the validation accuracy?

:SGD 71ps

```
Epoch [1/10], Step [71/71], Loss: 2.1925
Accuracy of the network on the 5000 validation images: 25.6 %
Epoch [2/10], Step [71/71], Loss: 1.8714
Accuracy of the network on the 5000 validation images: 38.0 %
Epoch [3/10], Step [71/71], Loss: 1.6779
Accuracy of the network on the 5000 validation images: 34.4 %
Epoch [4/10], Step [71/71], Loss: 1.5494
Accuracy of the network on the 5000 validation images: 39.0 %
Epoch [5/10], Step [71/71], Loss: 1.4755
Accuracy of the network on the 5000 validation images: 43.4 %
Epoch [6/10], Step [71/71], Loss: 1.4065
Accuracy of the network on the 5000 validation images: 47.2 %
Epoch [7/10], Step [71/71], Loss: 1.3365
Accuracy of the network on the 5000 validation images: 43.4 %
Epoch [8/10], Step [71/71], Loss: 1.2709
Accuracy of the network on the 5000 validation images: 48.6 %
Epoch [9/10], Step [71/71], Loss: 1.2077
Accuracy of the network on the 5000 validation images: 47.0 %
Epoch [10/10], Step [71/71], Loss: 1.1507
Accuracy of the network on the 5000 validation images: 50.8 %
```

:Adam 71ps

```
Epoch [1/10], Step [71/71], Loss: 1.8775
Accuracy of the network on the 5000 validation images: 38.6 %
Epoch [2/10], Step [71/71], Loss: 1.5067
Accuracy of the network on the 5000 validation images: 48.6 %
Epoch [3/10], Step [71/71], Loss: 1.3428
Accuracy of the network on the 5000 validation images: 49.8 %
Epoch [4/10], Step [71/71], Loss: 1.2231
Accuracy of the network on the 5000 validation images: 53.2 %
Epoch [5/10], Step [71/71], Loss: 1.0836
Accuracy of the network on the 5000 validation images: 59.6 %
Epoch [6/10], Step [71/71], Loss: 0.9696
Accuracy of the network on the 5000 validation images: 58.8 %
Epoch [7/10], Step [71/71], Loss: 0.8907
Accuracy of the network on the 5000 validation images: 63.6 %
Epoch [8/10], Step [71/71], Loss: 0.7791
Accuracy of the network on the 5000 validation images: 66.6 %
Epoch [9/10], Step [71/71], Loss: 0.7348
Accuracy of the network on the 5000 validation images: 62.4 %
Epoch [10/10], Step [71/71], Loss: 0.6338
Accuracy of the network on the 5000 validation images: 63.2 %
```

לעומת רכיב loss מתקדם נזק ור'  
לעומת loss סטנדרטי הנקרא SGD  
Adam ב 63.2 נטפס 50.8%.

- (b) (5 points) Again using Adam, train the network with a learning rate of 0.005 instead of the default 0.00005. How is the training process affected?

```
Epoch [1/10], Step [71/71], Loss: 22.3386
Accuracy of the network on the 5000 validation images: 13.0 %
Epoch [2/10], Step [71/71], Loss: 2.3835
Accuracy of the network on the 5000 validation images: 13.8 %
Epoch [3/10], Step [71/71], Loss: 2.3339
Accuracy of the network on the 5000 validation images: 16.6 %
Epoch [4/10], Step [71/71], Loss: 2.3180
Accuracy of the network on the 5000 validation images: 20.8 %
Epoch [5/10], Step [71/71], Loss: 2.2471
Accuracy of the network on the 5000 validation images: 20.2 %
Epoch [6/10], Step [71/71], Loss: 2.1806
Accuracy of the network on the 5000 validation images: 21.8 %
Epoch [7/10], Step [71/71], Loss: 2.1958
Accuracy of the network on the 5000 validation images: 14.4 %
Epoch [8/10], Step [71/71], Loss: 2.1598
Accuracy of the network on the 5000 validation images: 18.8 %
Epoch [9/10], Step [71/71], Loss: 2.1917
Accuracy of the network on the 5000 validation images: 16.2 %
Epoch [10/10], Step [71/71], Loss: 2.1279
Accuracy of the network on the 5000 validation images: 20.2 %
```

ויאגדה כי loss מתקדם יתבצע איזה מושג קשור לזרם  
נזק וחכמת ה-NN רצוי או לאו. מושג זה יתבצע  
וירגע.

(c) (5 points) Train the network using Adam and a step size of 0.00005, but change the architecture so that it doesn't include Batch Normalization or Dropout layers (simply comment out the relevant lines in the network's definition). How does the training loss evolve when compared to the one with BatchNorm + Dropout? Do these additions seem to accelerate the learning process?

```
Epoch [1/10], Step [71/71], Loss: 2.0798
Accuracy of the network on the 5000 validation images: 30.8 %
Epoch [2/10], Step [71/71], Loss: 1.7767
Accuracy of the network on the 5000 validation images: 38.8 %
Epoch [3/10], Step [71/71], Loss: 1.6116
Accuracy of the network on the 5000 validation images: 42.2 %
Epoch [4/10], Step [71/71], Loss: 1.5089
Accuracy of the network on the 5000 validation images: 49.2 %
Epoch [5/10], Step [71/71], Loss: 1.4111
Accuracy of the network on the 5000 validation images: 51.6 %
Epoch [6/10], Step [71/71], Loss: 1.3259
Accuracy of the network on the 5000 validation images: 49.0 %
Epoch [7/10], Step [71/71], Loss: 1.2276
Accuracy of the network on the 5000 validation images: 53.4 %
Epoch [8/10], Step [71/71], Loss: 1.1299
Accuracy of the network on the 5000 validation images: 54.8 %
Epoch [9/10], Step [71/71], Loss: 1.0410
Accuracy of the network on the 5000 validation images: 58.6 %
Epoch [10/10], Step [71/71], Loss: 0.9175
Accuracy of the network on the 5000 validation images: 54.8 %
```

השאלה מבקשת למדוד איזה השוני ב-<sup>epoch</sup> loss בין ה-BatchNorm + Dropout לבין ה-<sup>epoch</sup> loss ללא ה-<sup>epoch</sup> loss.

השאלה מבקשת למדוד איזה השוני ב-<sup>epoch</sup> loss בין ה-BatchNorm + Dropout לבין ה-<sup>epoch</sup> loss ללא ה-<sup>epoch</sup> loss.

סמןpdf