

1) Direct construction of Valid Kernels

$$\textcircled{a} \quad K(x, z) = K_1(x, z) + K_2(x, z)$$

\Rightarrow Proof:

Using Mercer's theorem,

K_1, K_2 are valid kernels, then their K_1 ,

K_2 (Gram matrixs) are positive-semidefinite

for all possible choices of the data set $\{x_i^{(n)}\}_{i=1}^N \{z_j^{(n)}\}_{j=1}^M$

$$a^T K_1 a = \sum_{i=1}^N \sum_{j=1}^M a_i K_{1,i,j} a_j \geq 0 \quad \forall a \in \mathbb{R}^N \quad \textcircled{1}$$

$$a^T K_2 a = \sum_{i=1}^N \sum_{j=1}^M a_i K_{2,i,j} a_j \geq 0 \quad \forall a \in \mathbb{R}^N \quad \textcircled{2}$$

adding $\textcircled{1}$ and $\textcircled{2}$

$$a^T (K_1 + K_2) a = \sum_{i=1}^N \sum_{j=1}^M a_i (K_{1,i,j} + K_{2,i,j}) a_j \geq 0 \quad \forall a \in \mathbb{R}^N$$

$$a^T (K) a = \sum_{i=1}^N \sum_{j=1}^M a_i (K_{i,j}) a_j \geq 0 \quad \forall a \in \mathbb{R}^N$$

Hence

$\underbrace{K(x, z)}_{\text{satisfy the Mercer's Theorem}} = K_1(x, z) + K_2(x, z)$ is a valid Kernel

Here Proved

$$(b) K(x, z) = K_1(x, z) - K_2(x, z)$$

Using ① and ② from part ①, subtracting ② from ①,

$$a^T(K_1 - K_2)a = \sum_{i=1}^N \sum_{j=1}^N a_i (\underbrace{K_{1,i,j} - K_{2,i,j}}_{\geq 0}) a_j, \quad \forall a \in \mathbb{R}^N$$

we cannot guarantee that $K_{1,i,j} - K_{2,i,j} \geq 0$

$\forall 1 \leq i \leq N$ and $1 \leq j \leq N$. Counter example, take

$$K_1 = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}, \quad K_2 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad \text{hence for these } K_1 \text{ and } K_2$$

$a^T(K_1 - K_2)a \leq 0$, semi-negative, Hence, it is not a valid Kernel

$$(c) K(x, x') = -a K_1(x, z)$$

Using eq ①,

$$b^T K_1 b = \sum_{i=1}^N \sum_{j=1}^N b_i K_{1,i,j} b_j \geq 0 \quad \forall b \in \mathbb{R}^N$$

Multiply the eq by $-a$, where $a \in \mathbb{R}^+$ be a positive real number.

$$b^T -a K_1 b = \sum_{i=1}^N \sum_{j=1}^N b_i -a K_{1,i,j} b_j \leq 0 \quad \forall b \in \mathbb{R}^N$$

Hence, it is semi-negative,

Hence it is not a valid Kernel function, because of $(-)$, but if $(-a)$ was a positive real number, then it would be a Kernel.

$$\textcircled{d} \quad K(x, z) = f(x) f(z)$$

\Rightarrow counter example,

$$f(x) = \sum_{i=1}^N x_i, \quad \forall x \in \mathbb{R}^N$$

$$\text{Let } x = \{1, -2, -1\}, \quad N=3, \quad z = \{1, 1, 1\}$$

using Mercer's theorem, the matrix K ,

$$K_{ij} = K(x^{(i)}, z^{(j)})$$

$$K = \begin{bmatrix} 1 & 1 & 1 \\ -2 & -2 & -2 \\ -1 & -1 & -1 \end{bmatrix} \quad \left. \begin{array}{l} \text{not symmetric} \\ \text{and not semi-positive definite} \end{array} \right\}$$

Hence, Kernel is not a valid Kernel.

$$\textcircled{e} \quad K(x, z) = K_3(\phi(x), \phi(z))$$

$$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$$

$\Rightarrow K_3$ is a valid Kernel,

$$K(x, z) \text{ is symmetric, i.e. } K_3(\phi(x), \phi(z)) = K_3(\phi(z), \phi(x))$$

This is true as K_3 is Kernel so it is symmetric.

Now, Positive-semi definite from Mercer's Theory,

$$a^T K a \geq 0 \quad \forall a \in \mathbb{R}^N, \text{ where}$$

$K_{ij} = \kappa_3(x^{(i)}, x^{(j)})$, K is symmetric and positive semi definite.

For, $\kappa(x, z) = \kappa_3(\phi(x), \phi(z))$,
 we construct K' , where $K'_{ij} = \kappa_3(\phi(x^{(i)}), \phi(z^{(i)}))$
 since κ_3 is valid kernel, then K' will be
 a semi-positive definite matrix, hence

$$a^T K' a \geq 0, \forall a \in \mathbb{R}^N.$$

(g) $\kappa(x, z) = \kappa_1(x, z) \kappa_2(x, z)$

$$\text{Let } K_1(x, z) = \langle \phi^1(x), \phi^1(z) \rangle \text{ and}$$

$$K_2(x, z) = \langle \phi^2(x), \phi^2(z) \rangle, \begin{matrix} x \in \mathbb{R}^m \\ z \in \mathbb{R}^n \end{matrix}$$

$$\begin{aligned} \kappa_1(x, z) \kappa_2(x, z) &= \langle \phi^1(x), \phi^1(z) \rangle \cdot \langle \phi^2(x), \phi^2(z) \rangle \\ &= \left(\sum_{i=1}^m \phi_i^1(x) \phi_i^1(z) \right) \left(\sum_{j=1}^n \phi_j^2(x) \phi_j^2(z) \right) \\ &= \sum_{i=1}^m \sum_{j=1}^n (\phi_i^1(x) \phi_j^2(x)) \cdot (\phi_i^1(z) \phi_j^2(z)) \\ &= \langle \phi(x), \phi(z) \rangle \end{aligned}$$

where, $\phi(x) = \begin{bmatrix} \phi_1^1(x) & \phi_1^2(x) \\ \phi_1^1(x) & \phi_2^2(x) \\ \phi_2^1(x) & \phi_1^2(x) \\ \phi_1^1(x) & \phi_3^2(x) \\ \phi_2^1(x) & \phi_2^2(x) \\ \phi_3^1(x) & \phi_1^2(x) \\ \vdots & \end{bmatrix}$

ϕ is a feature map of Δ^{2n+1} , where

$$\phi_{ij}(x) = \phi_i^1(x) \phi_j^2(x), \quad 1 \leq i, j \leq n$$

Hence, $K(x, z)$ is a Kernel for a feature map $\phi(\cdot)$.

(\Leftarrow) $K(x, z) = P(K_1(x, z))$

$P: \mathbb{R} \rightarrow \mathbb{R}$ with positive coefficients.

\Rightarrow

since, polynomial term is a product of kernels with a positive kernel and then sum, i.e

$$P(K_1(x, z)) = a_n(K_1(x, z))^n + a_{n-1}(K_1(x, z))^{n-1} + \dots$$

$\dots + a_0$, where

$a_n, a_{n-1}, \dots, a_0 > 0$ are positive coefficient

So, using property of ⑨ i.e

$\kappa(x, z) = \kappa_1(x, z) \kappa_2(x, z)$ is a Kernel

and property of ⑩ $\kappa(x, z) = a \kappa_1(x, z)$

is a Kernel, where a is positive real number

and property of ⑪ $\kappa(x, z) = \kappa_1(x, z) + \kappa_2(x, z)$

, we can say

$\delta(\kappa_1(x, z))$ is also a Kernel

Hence Proved

h) $\kappa(x, z) = (x^T z + 1)^2$ is a Kernel,

feature map ϕ associated with $\kappa(x, z)$ s.t

$$\kappa(x, z) = \phi(x)^T \phi(z), \text{ let } D=2$$

$$\Rightarrow x = \{x_1, x_2\}, z = \{z_1, z_2\}$$

$$\begin{aligned}\kappa(x, z) &= (x_1 z_1 + x_2 z_2 + 1)^2 \\ &= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 + \\ &\quad 2x_1 z_1 + 2x_2 z_2 + 1 \\ &= \phi(x)^T \phi(z).\end{aligned}$$

where $\phi(x) = (x_1^2, x_2^2, \sqrt{2}x_1 x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$ is a feature map.

(i) Extra Credit

$$k(x, z) = \exp(-\|x - z\|^2 / 2\sigma^2)$$

$= \phi(x)^T \phi(z)$, find the feature vector $\phi(x)$

\Rightarrow using Taylor series, $\exp(x) = \sum_{n=0}^{\infty} \frac{1}{n!} x^n$

$$k(x, z) = \exp\left(\frac{-x^T x - z^T z + 2x^T z}{2\sigma^2}\right).$$

$$= \exp\left(\frac{-x^T x}{2\sigma^2}\right) \exp\left(\frac{-z^T z}{2\sigma^2}\right) \exp\left(\frac{x^T z}{\sigma^2}\right).$$

$$= \exp\left(\frac{-x^T x}{2\sigma^2}\right) \exp\left(\frac{-z^T z}{2\sigma^2}\right) \left(\sum_{n=0}^{\infty} \frac{1}{n!} \left(\frac{x^T z}{\sigma^2}\right)^n \right)$$

so

$$\phi(x) = \exp\left(\frac{-x^T x}{2\sigma^2}\right) A, \text{ where } A \text{ is}$$

$$A = \left[1, \dots, \sqrt{\frac{1}{P! \sigma^{2P}}} \begin{pmatrix} P \\ j_1, j_2, \dots, j_m \end{pmatrix} x_1^{j_1} x_2^{j_2} \dots x_m^{j_m}, \dots \right]$$

2) Implementing soft margin SVM by optimizing Priority Objective

(a) Derivatives of loss function $E(\omega, b)$ w.r.t ω, b .

$$E(\omega, b) = \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \max(0, 1 - y^{(i)}(\omega^\top \phi(x^{(i)}) + b))$$

(i) Taking derivative w.r.t "w"

$$\Rightarrow \nabla_\omega E(\omega, b) = \nabla_\omega \left(\frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \max(0, 1 - y^{(i)}(\omega^\top \phi(x^{(i)}) + b)) \right)$$

$$\Rightarrow \nabla_\omega E(\omega, b) = \omega + C \sum_{i=1}^n \frac{d}{d\omega} (\max(0, 1 - y^{(i)}(\omega^\top \phi(x^{(i)}) + b)))$$

$$\nabla_\omega E(\omega, b) = \omega - C \sum_{i=1}^n \mathbb{I}[1 - y^{(i)}(\omega^\top \phi(x^{(i)}) + b) > 0] y^{(i)} \phi(x^{(i)})$$

$$\nabla_\omega E(\omega, b) = \omega - C \sum_{i=1}^n \mathbb{I}[y^{(i)}(\omega^\top \phi(x^{(i)}) + b) > 1] y^{(i)} \phi(x^{(i)})$$

(ii) Taking derivative w.r.t "b"

$$\frac{\partial E(w, b)}{\partial b} = \frac{\partial}{\partial b} \left(\frac{1}{2} \|w\|^2 \right) + C \sum_{i=1}^n \frac{\partial}{\partial b} \left(\max(0, 1 - y^{(i)}(w^T \phi(x^{(i)}) + b)) \right)$$

$$= -C \sum_{i=1}^n \mathbb{I} \left[1 - y^{(i)}(w^T \phi(x^{(i)}) + b) > 0 \right] y^{(i)}$$

$$\frac{\partial E(w, b)}{\partial b} = -C \sum_{i=1}^n \mathbb{I} \left[y^{(i)}(w^T \phi(x^{(i)}) + b) < 1 \right] y^{(i)}$$

(b) Autograder

(c)

[NumEpochs: 1] Accuracy: 54.17%

b: [0.01], W: [[0.224 -0.0855 0.545 0.206]]

[NumEpochs: 3] Accuracy: 54.17%

b: [0.02], W: [[0.44848122 -0.17019759 1.08918105 0.41163421]]

[NumEpochs: 10] Accuracy: 95.83%

b: [-0.14], W: [[-0.1648026 -0.80606447 1.37816462 0.57445096]]

[NumEpochs: 30] Accuracy: 95.83%

b: [-0.175], W: [[-0.20885861 -0.69979483 1.30489009 0.56605151]]

[NumEpochs: 100] Accuracy: 95.83%

b: [-0.315], W: [[-0.28240917 -0.77529188 1.75856715 0.82441652]]

3)

Asymmetric SVM

$$\min_{w, b, \xi} \frac{1}{2} w^T w + C_0 \sum_{i: y^{(i)} = -1} \xi^{(i)} + C_1 \sum_{i: y^{(i)} = 1} \xi^{(i)}$$

$$\text{s.t. } y^{(i)}(w^T \phi(x^{(i)}) + b) \geq 1 - \xi^{(i)}, \forall i=1, \dots, N$$

$$\xi^{(i)} \geq 0 \quad \forall i=1, \dots, N$$

a)

Defining/

$$g_i(x) = 1 - \xi^{(i)} - y^{(i)}(w^T \phi(x^{(i)}) + b)$$

$$h_i(x) = -\xi^{(i)}$$

so,

$$\mathcal{L}(w, b, \xi, \alpha, \mu) = \frac{1}{2} w^T w + C_0 \sum_{i: y^{(i)} = -1} \xi^{(i)} + C_1 \sum_{i: y^{(i)} = 1} \xi^{(i)} + \sum_{i=1}^n \alpha^{(i)} g_i(x) + \sum_{i=1}^n \mu^{(i)} h_i(x)$$

here $\alpha^{(i)}, \mu^{(i)} \geq 0 \quad \forall i = \{1, \dots, N\}$

b)

$$\nabla_w \mathcal{L}(w, b, \xi, \alpha, \mu) = w + \sum_{i=1}^n \alpha^{(i)} \frac{\partial g_i(x)}{\partial w}$$

$$\nabla_w \mathcal{L}(w, b, \xi, \alpha, \mu) = w - \sum_{i=1}^n \alpha^{(i)} (y^{(i)} \phi(x^{(i)}))$$

ii)

$$\frac{\partial \mathcal{L}(\omega, b, \xi, \alpha, \mu)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} \frac{d g_i(x)}{db}$$

$$\frac{\partial \mathcal{L}(\omega, b, \xi, \alpha, \mu)}{\partial b} = -\sum_{i=1}^N \alpha^{(i)} y^{(i)}$$

iii) $\nabla_{\xi^{(i)}} \mathcal{L}(\omega, b, \xi, \alpha, \mu)$

$\nabla_{\xi^{(i)}} \mathcal{L}(\omega, b, \xi, \alpha, \mu) =$

$$\sum_{i=1}^N c_0 \mathbb{I}[y^{(i)} = -1] + c_1 \mathbb{I}[y^{(i)} = 1] - \sum_{i=1}^N \alpha^{(i)} - \sum_{i=1}^N \mu^{(i)}$$

c) Dual optimization problem,

$$P^* = \max_{\substack{\alpha, \mu \\ \alpha \geq 0 \\ \mu \geq 0}} \min_{w, b} \left(\frac{1}{2} w^T w + c_0 \sum_{i: y^{(i)} = -1} \xi^{(i)} + c_1 \sum_{i: y^{(i)} = 1} \xi^{(i)} + \sum_{i=1}^N \alpha^{(i)} g_i(x) + \sum_{i=1}^N \mu^{(i)} h_i(x) \right)$$

\Rightarrow solving min first, $\nabla_w \mathcal{L}(\omega, b, \xi, \alpha, \mu) = 0$ and
 $\frac{\partial \mathcal{L}(\omega, b, \xi, \alpha, \mu)}{\partial b} = 0$

$$w = \sum_{i=1}^n \alpha^{(i)} y^{(i)} \phi(x^{(i)}) \quad \text{and} \quad b = \sum_{i=1}^n \alpha^{(i)} y^{(i)}$$

$$\Rightarrow \max_{\substack{\alpha, \mu \\ \alpha^{(i)} \geq 0 \\ \mu^{(i)} \geq 0 \\ \#i}} \tilde{\mathcal{L}}(\alpha) = \sum_{i=1}^n \alpha^{(i)} (1 - \xi_i^{(i)}) + \sum_{i=1}^n \mu^{(i)} (-\xi_i^{(i)}) + \sum_{i=1}^n (c_0 \mathbb{I}[y^{(i)} = -1] + c_1 \mathbb{I}[y^{(i)} = 1]) \xi_i^{(i)} - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \phi(x^{(i)})^T \phi(x^{(j)})$$

This can be rewritten as below by removing redundant variable μ .

$$\max_{\substack{\alpha \\ \alpha^{(i)} \geq 0 \\ \#i}} \tilde{\mathcal{L}}(\alpha) = \sum_{i=1}^n \alpha^{(i)} (1 - 2\xi_i^{(i)}) + \sum_{i=1}^n (c_0 \mathbb{I}[y^{(i)} = -1] + c_1 \mathbb{I}[y^{(i)} = 1]) \xi_i^{(i)} - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} \phi(x^{(i)})^T \phi(x^{(j)})$$

Y)

SVMs with Convex Optimization

(a)

$$\underset{\varphi}{\text{minimize}} \quad \frac{1}{2} \varphi^T P \varphi + q^T \varphi$$

$$\text{subject to} \quad q \varphi \leq h$$

$$A \varphi = b$$

$$\underset{\alpha}{\text{minimize}} \quad \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha^{(n)} \alpha^{(m)} y^{(n)} y^{(m)} k(x^{(n)}, x^{(m)}) - \sum_{n=1}^N \alpha^{(n)}$$

$$\text{subject to} \quad 0 \leq \alpha^{(n)} \leq C$$

$$\sum_{n=1}^N \alpha^{(n)} y^{(n)} = 0$$

Here, $b = 0$, $v = [\alpha^1, \alpha^2, \dots, \alpha^N]^T$ ($N \times 1$) matrix

$$A = [y^1, y^2, \dots, y^N]^T \quad (1 \times N) \text{ matrix}$$

$$q_v = [-1, 1, \dots, -1]^T \quad (N \times 1) \text{ matrix}$$

$$h = [0, 0, \dots, 0, c, -c, \dots, -c]^T \quad ((N+N) \times 1) \text{ matrix}$$

$$Q = \begin{bmatrix} -I_{N \times N} \\ I_{N \times N} \end{bmatrix} \quad ((N+N) \times N) \text{ matrix; } I \text{ is a identity matrix}$$

$$P = \begin{bmatrix} p_{1,1} & \dots & p_{1,N} \\ \vdots & & \vdots \\ p_{N,1} & \dots & p_{N,N} \end{bmatrix}, \quad (N \times N \text{ matrix})$$

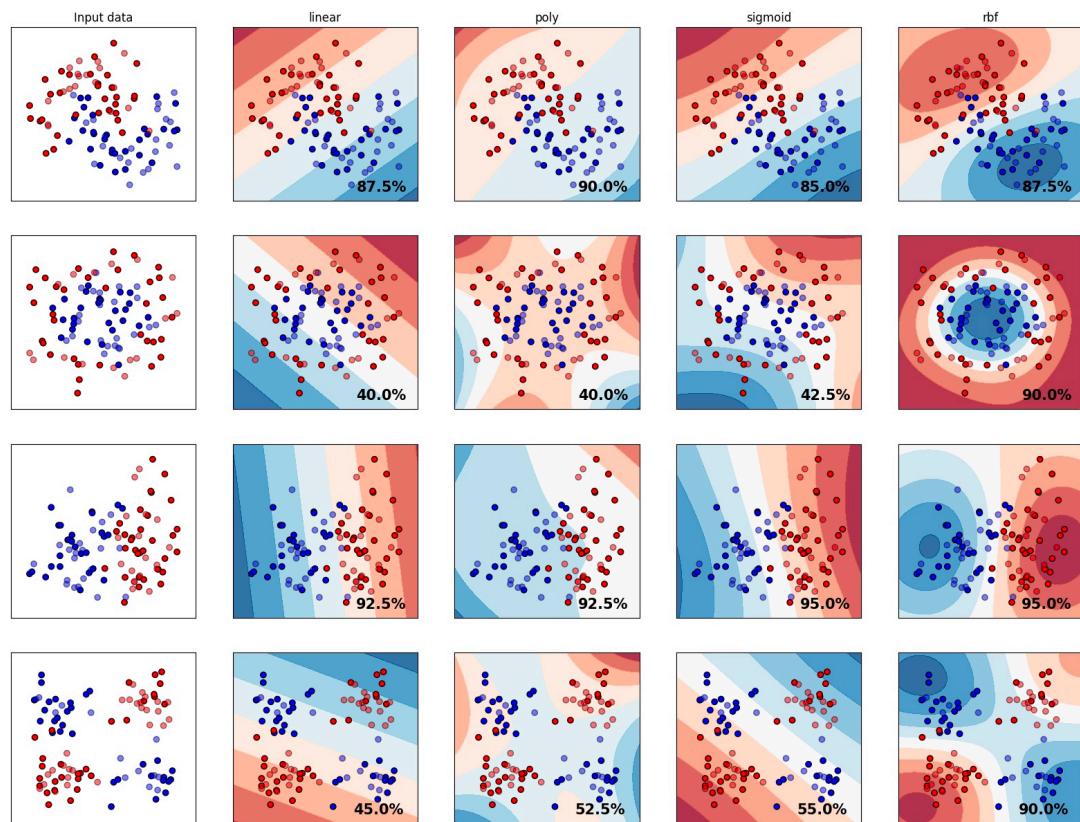
$$\text{where } p_{i,j} = y^{(i)} y^{(j)} k(x^{(i)}, x^{(j)})$$

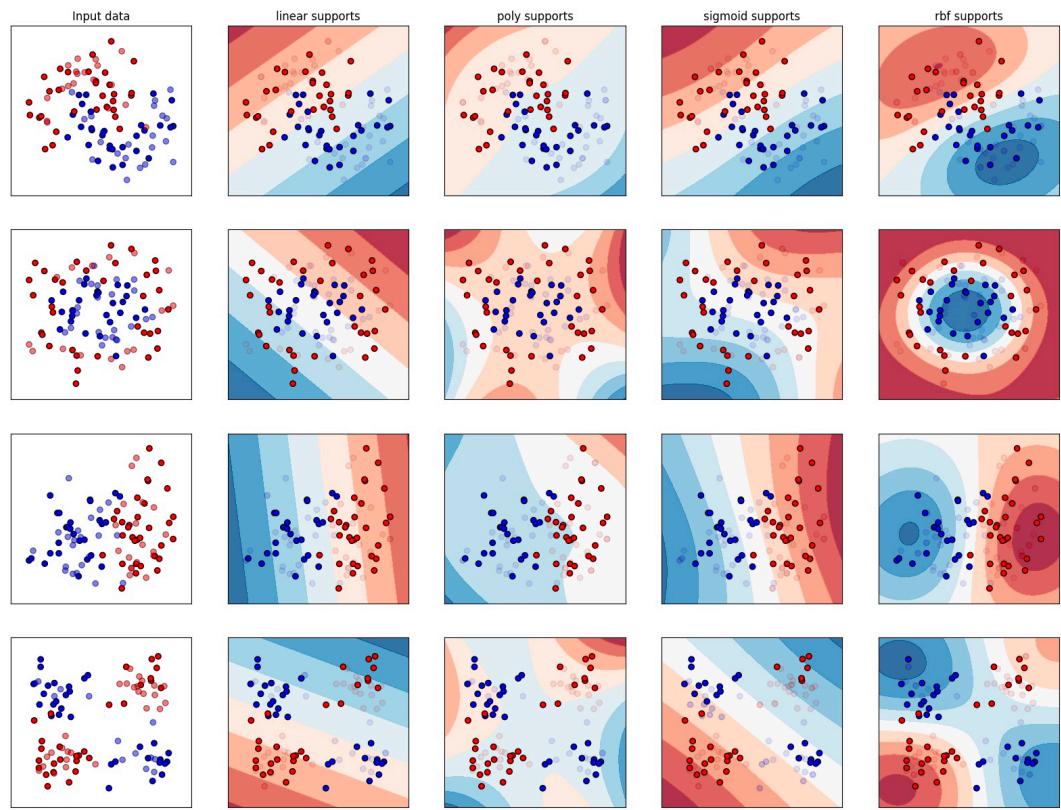
(b)

Auto grader

(c)

	linear	poly	sigmoid	rbf
dataset 0	87.50	90.00	85.00	87.50
dataset 1	40.00	40.00	42.50	90.00
dataset 2	92.50	92.50	95.00	95.00
dataset 3	45.00	52.50	55.00	90.00





5) Neural Network layer Implementation

(a)

$$Y = XW + B \quad , \quad Y \in \mathbb{R}^{N \times D_{\text{out}}}$$

$$B \in \mathbb{R}^{N \times D_{\text{out}}}$$

$$X \in \mathbb{R}^{N \times D_{\text{in}}}$$

$$W \in \mathbb{R}^{D_{\text{in}} \times D_{\text{out}}}$$

$L \in \mathbb{R}$ loss function

Let $\frac{\partial L}{\partial W_{i,j}}$ is the i^{th} row, j^{th} column element of $\frac{\partial L}{\partial W} \in \mathbb{R}^{D_{\text{in}} \times D_{\text{out}}}$

$$1 \leq n \leq N \quad , \quad 1 \leq i \leq D_{\text{in}} \quad , \quad 1 \leq j \leq D_{\text{out}}$$

$$\Rightarrow \frac{\partial L}{\partial W_{i,j}} = \sum_{n=1}^N \frac{\partial L}{\partial Y_j^{(n)}} \frac{\partial Y_i^{(n)}}{\partial W_{i,j}} \quad , \quad \text{here}$$

$$\text{where } Y_j^{(n)} = X^{(n)} W_j + B_j^{(n)}$$

$Y_j^{(n)}$ is the output of j^{th} neuron in the n^{th} sample in a batch.

$$\frac{\partial L}{\partial W_{i,j}} = \sum_{n=1}^N \frac{\partial L}{\partial Y_j^{(n)}} X_i^{(n)} = X_i^T \frac{\partial L}{\partial Y_j}$$

Hence

$$\boxed{\frac{\partial L}{\partial W} = X^T \frac{\partial L}{\partial Y}}$$

$$\frac{\partial \mathcal{L}}{\partial b_j}$$

$$\frac{\partial \mathcal{L}}{\partial b_j} = \sum_{n=1}^N \sum_{m=1}^{D_{\text{out}}} \frac{\partial \mathcal{L}}{\partial Y_m^{(n)}} \frac{\partial Y_m^{(n)}}{\partial b_j}$$

where $Y_m^{(n)} = X^{(n)} W_m + B_m^{(n)}$ | $\frac{\partial Y_m^{(n)}}{\partial b_j} = 1$

$$\frac{\partial \mathcal{L}}{\partial b_j} = \sum_{n=1}^N \sum_{m=1}^{D_{\text{out}}} \frac{\partial \mathcal{L}}{\partial Y_m^{(n)}} \times 1$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial B} = \sum_{n=1}^N \frac{\partial \mathcal{L}}{\partial Y^{(n)}}}$$

$$\frac{\partial \mathcal{L}}{\partial X_i}$$

$$\frac{\partial \mathcal{L}}{\partial X_i^n} = \sum_{n=1}^N \sum_{m=1}^{D_{\text{out}}} \frac{\partial \mathcal{L}}{\partial Y_m^{(n)}} \frac{\partial Y_m^{(n)}}{\partial X_i^n}$$

where $Y_m^{(n)} = X^{(n)} W_m + B_m^{(n)}$

as the n^{th} -sample in X is only related with the n^{th} -sample in Y and $\frac{\partial Y_m^{(n)}}{\partial X_i^n} = 0 \forall n \neq n$ so

$$\frac{\partial \mathcal{L}}{\partial x_i^{(m)}} = \sum_{m=1}^{\text{Dout}} \frac{\partial \mathcal{L}}{\partial y_m^{(m)}} \frac{\partial y_m^{(m)}}{\partial x_i^{(m)}}$$

$$\frac{\partial \mathcal{L}}{\partial x_i^{(m)}} = \sum_{m=1}^{\text{Dout}} \frac{\partial \mathcal{L}}{\partial y_m^{(m)}} (w_{i,m})$$

$$\frac{\partial \mathcal{L}}{\partial x_i^{(m)}} = \frac{\partial \mathcal{L}}{\partial y^{(m)}} (w_i)^T$$

$$\frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} w^T$$

(b) $x \in \mathbb{R}^{N \times D}$, $y = \text{ReLU}(x)$
 $x \in X$, $y = \text{ReLU}(x) = \max(0, x)$

$$\frac{\partial \mathcal{L}}{\partial x} \in \mathbb{R}^{N \times D}, \quad 1 \leq i \leq N, 1 \leq j \leq D$$

Let $\frac{\partial \mathcal{L}}{\partial x_{ij}} = \left(\frac{\partial \mathcal{L}}{\partial y_{ij}} \right) \frac{\partial y_{ij}}{\partial x_{ij}}$ } True as ReLU is applied elementwise,

$$\frac{\partial \mathcal{L}}{\partial x_{ij}} = \frac{\partial \mathcal{L}}{\partial y_{ij}} \cdot \mathbb{I}(x_{ij} > 0)$$

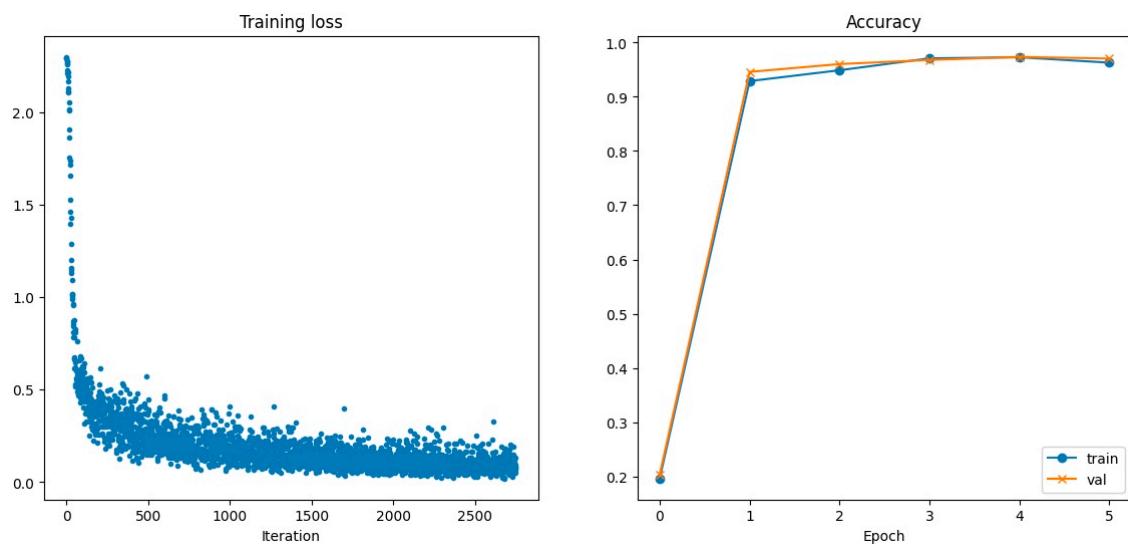
\odot is element wise multiplication

$$\frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} \odot \mathbb{I}(x > 0)$$

(c)

Autograder

(d)



(e)

(Epoch 0 / 5) train acc: 19.60% val_acc: 20.36%

(Epoch 1 / 5) train acc: 92.90% val_acc: 94.60%

(Epoch 2 / 5) train acc: 94.90% val_acc: 96.04%

(Epoch 3 / 5) train acc: 97.10% val_acc: 96.80%

(Epoch 4 / 5) train acc: 97.30% val_acc: 97.36%

(Epoch 5 / 5) train acc: 96.30% val_acc: 97.06%

Accuracy
on
training set

Test accuracy: 96.31%

Accuracy on
test set.