

Q1

 $L_1(a)$

$$L_P = \frac{1}{2} W^T S W - \nu \rho + \frac{1}{N} \sum_t \xi^t$$

$$\text{subj to } \begin{cases} r^t (w^T x^t + w_0) \geq \rho - \xi^t \\ \xi^t \geq 0 \\ \rho \geq 0, \nu \in [0, 1] \end{cases}$$

$$L_P = \frac{1}{2} W^T S W - \nu \rho + \frac{1}{N} \sum_t \xi^t - \sum_t \alpha^t [r^t (w^T x^t + w_0) - \rho + \xi^t] - \sum_t \mu^t \xi^t - \eta \rho$$

$\because S$ is symmetric positive definite

$$\left\{ \begin{aligned} \frac{\partial L_P}{\partial w} &= S w - \sum_t \alpha^t r^t x^t = 0 \Rightarrow S w = \sum_t \alpha^t r^t x^t \end{aligned} \right.$$

$$\left\{ \begin{aligned} \frac{\partial L_P}{\partial w_0} &= -\sum_t \alpha^t r^t = 0 \Rightarrow \sum_t \alpha^t r^t = 0 \quad \Rightarrow \text{cancel } w_0 \end{aligned} \right.$$

$$\left\{ \begin{aligned} \frac{\partial L_P}{\partial \rho} &= -\nu - \sum_t \alpha^t (-1) - \eta = 0 \Rightarrow \nu = \sum_t \alpha^t - \eta \Rightarrow \text{cancel } \rho \end{aligned} \right.$$

$$\left\{ \begin{aligned} \frac{\partial L_P}{\partial \xi^t} &= \frac{1}{N} - \alpha^t - \mu^t = 0 \Rightarrow \mu^t = \frac{1}{N} - \alpha^t \quad \Rightarrow \text{cancel } \xi^t \end{aligned} \right.$$

$$L_P = \frac{1}{2} W^T S W - \left(\sum_t \alpha^t - \eta \right) \rho + \frac{1}{N} \sum_t \xi^t - \sum_t \alpha^t [r^t (w^T x^t + w_0) - \rho + \xi^t] - \sum_t \left(\frac{1}{N} - \alpha^t \right) \xi^t$$

$$= \frac{1}{2} W^T S W - \sum_t \alpha^t r^t w^T x^t$$

$$= \frac{1}{2} W^T \sum_t \alpha^t r^t x^t - \sum_t \alpha^t r^t w^T x^t$$

$$= -\frac{1}{2} \sum_t \alpha^t r^t w^T x^t$$

$$= -\frac{1}{2} \sum_t \alpha^t r^t w^T S^{-1} S x^t$$

$$= -\frac{1}{2} \sum_t \alpha^t r^t (S w)^T S^{-1} x^t$$

$$= -\frac{1}{2} \sum_t \alpha^t r^t \left[\sum_k \alpha^k r^k (x^k)^T \right] S^{-1} x^t$$

$$L_P = -\frac{1}{2} \sum_t \sum_k \alpha^t \alpha^k r^t r^k (x^k)^T S^{-1} x^t$$

$$\text{subj. } \sum_t \alpha^t r^t = 0, 0 \leq \alpha^t \leq \frac{1}{N}, \sum_t \alpha^t \geq \nu$$

Q2.

(b)

$$L_p = \frac{1}{2} w^T S w + \sum_t c^t z^t - \sum_t \alpha^t [r^t (w^T x^t + w_0) - 1 + z^t] - \sum_t \mu^t z^t$$

$$\text{subj. } z^t \geq 0, r^t (w^T x^t + w_0) \geq 1 - z^t$$

$\therefore S$ is positive definite.

$$\frac{\partial L_p}{\partial w} = S w - \sum_t \alpha^t r^t x^t = 0 \Rightarrow S w = \sum_t \alpha^t r^t x^t$$

$$\frac{\partial L_p}{\partial w_0} = - \sum_t \alpha^t r^t = 0 \Rightarrow \sum_t \alpha^t r^t = 0 \Rightarrow \text{cancel } w_0$$

$$\frac{\partial L_p}{\partial z^t} = c^t - \alpha^t - \mu^t = 0 \Rightarrow \mu^t = c^t - \alpha^t \Rightarrow \text{cancel } z^t$$

$$L_p = \frac{1}{2} w^T S w + \sum_t c^t z^t - \sum_t \alpha^t [r^t (w^T x^t + w_0) - 1 + z^t] - \sum_t (c^t - \alpha^t) z^t$$

$$= \frac{1}{2} w^T \sum_t \alpha^t r^t x^t - \sum_t \alpha^t r^t w^T x^t + \sum_t \alpha^t$$

$$= \sum_t \alpha^t - \frac{1}{2} w^T \sum_t \alpha^t r^t x^t$$

$$w^T = \sum_t \alpha^t r^t (x^t)^T S^{-1}$$

$$L_p = \sum_t \alpha^t - \frac{1}{2} \sum_k \alpha^k r^k (x^k)^T S^{-1} \sum_t \alpha^t r^t x^t$$

$$= \sum_t \alpha^t - \frac{1}{2} \sum_k \sum_t \alpha^k \alpha^t r^k r^t (x^k)^T S^{-1} x^t$$

$$\text{subj. } \sum_t \alpha^t r^t = 0, 0 \leq \alpha^t \leq c^t$$

For perceptron:
the update rule is.

$$\begin{cases} w = 0 \\ f(x) = \langle w, x \rangle + b \\ \text{for } k \text{ epoch} = 1, 2, \dots, T \\ \text{for all } (x_i, r_i), \text{ where } i = 1, \dots, n \\ \text{if } f(x_i) \cdot r_i < 0 \\ w \leftarrow w + r_i x_{ik} \end{cases}$$

$$w = 0 + \alpha_1 r_1 x_1 + \alpha_2 r_2 x_2 + \dots + \alpha_n r_n x_n \\ = \sum_{i=1}^n \alpha_i r_i x_i$$

α_i is number of times perceptron made mistakes on i th observation within entirely T loops.

New update rule with respect to α_i

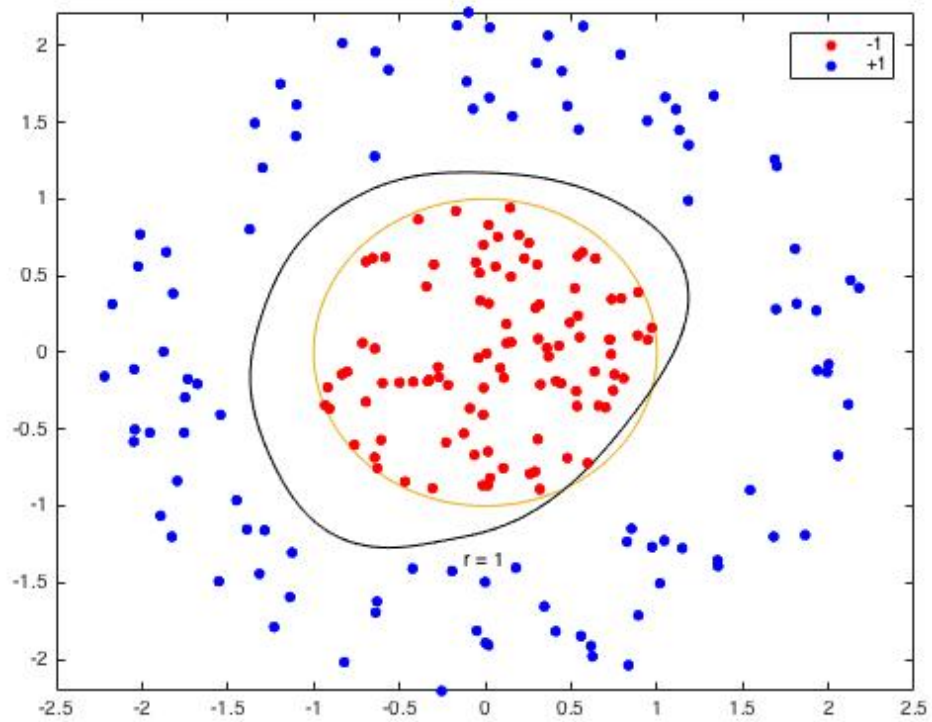
$$\begin{cases} \alpha = 0 \\ f(x) = \sum_j \alpha_j r_j \langle x_j, x \rangle + b \\ \text{for all } (x_i, r_i), \text{ where } i = 1, \dots, n \\ \text{if } f(x_i) \cdot r_i < 0 \\ \alpha_i \leftarrow \alpha_i + 1 \end{cases}$$

replace

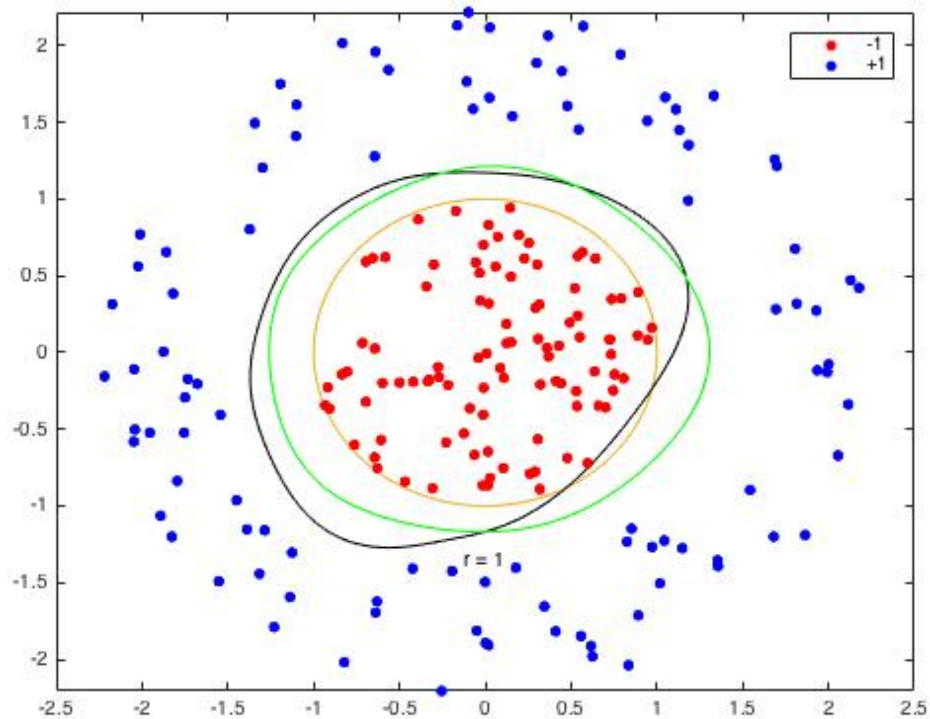
If we replace inner product with kernel function $\phi(x_j, x)$, so, $f(x) = b + \sum_j \alpha_j r_j \phi(x_j, x)$

Q3

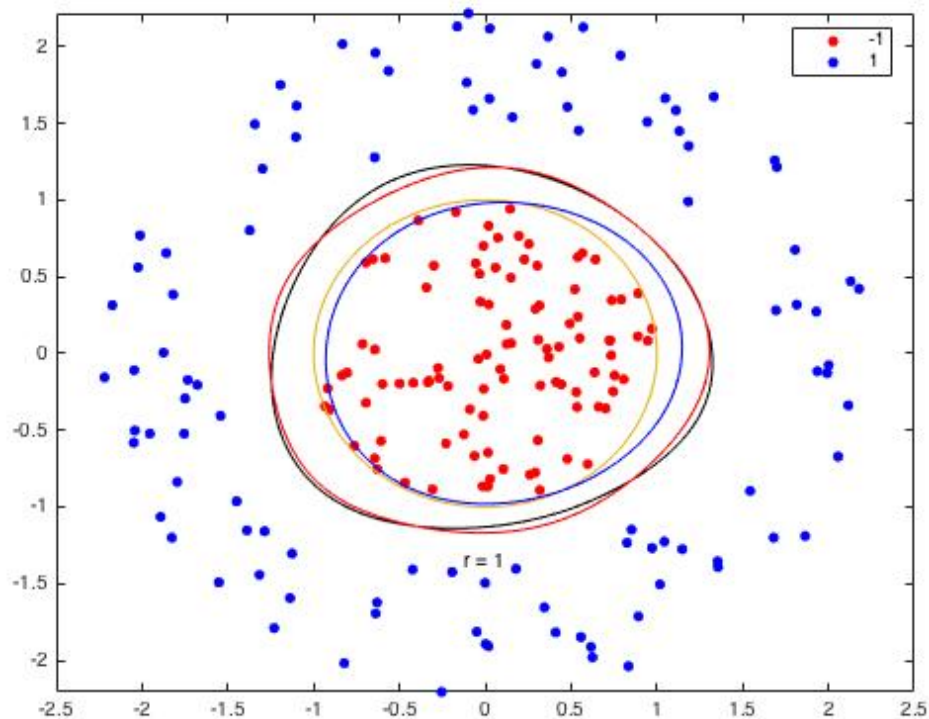
For data3, the error rate =0. The decision boundary for GD kernel is the black curve:



The decision boundary for built-in SVM is the green curve. The decision boundary for GD is the black curve. We can discover that margin around the green curve(SVM) is larger than black curve(perceptron).



When I play with the box constraint parameter, the smaller the box constraint is, the larger the boundary is. The blue curve represents box constraint=0.01 and black represents 1 while the red 100. The box constraint means the value of C in object function $\frac{1}{2} \langle w, w \rangle + C \sum_i s_i$, where C is penalty we choose for slack variable. If box constraint or C is large, it means we penalize the severely for slack and we can observe less margin-violating observations and less support vectors. So you can see it as soft margin.



In addition,

On the digits49_train data, the error rate= 0.0047

On the digits49_test data, the error rate= 0.0282

On the digits79_train data, the error rate= 0.0035

On the digits79_test data, the error rate= 0.0035