



$$2.1 \quad f(x) = \max_i f_i(x), \quad i=1, \dots, m.$$

$$f(x) \geq f_k(x), \quad k=1, \dots, m.$$

Suppose that at a point  $x_0$ ,  $f(x_0) = f_k(x_0)$  and  $\sqrt{\text{one of the}}$  subgradients of  $f_k(x)$  at  $x_0$  as  $g$ .

$$\text{We have: } f(x) \geq f_k(x) \geq f_k(x_0) + g^T(x - x_0) = f(x_0) + g^T(x - x_0)$$

So  $g$  also satisfies the definition of subgradient for  $f(x)$ .

$$2.2 \quad J(w) = \max \{0, 1 - yw^T x\}.$$

$$\text{When } 1 - yw^T x \leq 0, J(w) = 0, \partial 0 = 0.$$

$$\text{When } 1 - yw^T x > 0, J(w) = 1 - yw^T x, \partial(1 - yw^T x) = -yx$$

$$g = \begin{cases} 0, & 1 - yw^T x \leq 0 \\ -yx, & 1 - yw^T x > 0 \end{cases}$$

can be one of  $J(w)$ 's subgradients

3.1 If  $\{x | w^T x = 0\}$  is a separating hyperplane for  $D$ .

then  $y_i w^T x_i \geq 0, i=1, \dots, n$ .

$$\ell(\hat{y}, y) = \max\{0, -\hat{y}y\} = \max\{0, -yw^T x\}. \text{ so } \ell(\hat{y}_i, y_i) = 0 \text{ and empirical}$$

loss is 0. Moreover,  $\ell(\hat{y}, y) \geq 0$  so every separating hyperplane is an empirical risk minimizer.



$$3.2 \ell(\hat{y}, y) = \max\{0, -\hat{y}y\} = \max\{0, -y w^T x\}$$

$$g = \begin{cases} 0, & y w^T x > 0 \\ -y x, & y w^T x < 0 \end{cases} \quad \text{can be one of the subgradients.}$$

Only need to update  $w$  to  $w + \eta y x_i$  when  $y_i w^T x_i < 0$ .

Choose a proper step size, we can implement stochastic subgradient descent now.

3.3. If  $w$  is initialized as 0,  $w$  should be linear combination of  $\{x_1, \dots, x_n\}$ . Because it's added a scaled  $x_i$  in each update.

$$6.1 J_i(w) = \frac{\lambda}{2} \|w\|^2 + \max\{0, 1 - y_i w^T x_i\}$$

$$\nabla J_i(w) = \begin{cases} \lambda w, & 1 - y_i w^T x_i < 0 \\ \lambda w - y_i x_i, & 1 - y_i w^T x_i > 0 \end{cases}, \text{ when } y_i w^T x_i = 1, \text{ it's undefined.}$$

$$6.2 \partial \left( \frac{\lambda}{2} \|w\|^2 \right) = \lambda w. \quad \partial \max\{0, 1 - y_i w^T x_i\} \ni \begin{cases} 0, & y_i w^T x_i \geq 1 \\ -y_i x_i, & \text{else} \end{cases}$$

$$g = \begin{cases} \lambda w, & y_i w^T x_i \geq 1 \\ \lambda w - y_i x_i, & \text{else} \end{cases} \in \partial J_i(w).$$





6.3 For step size  $\eta_t = \frac{1}{\lambda t}$ ,  $w^{(t+1)} = w^{(t)} - \eta_t g$

$$\eta_t g = \begin{cases} \frac{1}{t} w^{(t)} - \frac{1}{\lambda t} y_i x_i, & y_i w^{(t)} x_i < 1 \\ \frac{1}{t} w^{(t)} & , \text{ else. } \end{cases}$$

$$w^{(t+1)} = \begin{cases} (1 - \frac{1}{t}) w^{(t)} + \frac{1}{\lambda t} y_i x_i, & y_i w^{(t)} x_i < 1 \\ (1 - \frac{1}{t}) w^{(t)} & , \text{ else. } \end{cases}$$

6.5.  $S_{t+1} = (1 - \eta_t \lambda) S_t$  and let  $w^{(t+1)} = S_t \cdot W_t$ .

When  $y_i w^{(t+1)} x_i \geq 1$ , only need to update  $S_t$ .

When  $y_i w^{(t+1)} x_i < 1$ ,  $S_{t+1} W_{t+1} = (1 - \eta_t \lambda) S_t \cdot W_t + \eta_t y_i x_i$ .

So  $W_{t+1} = W_t + \frac{1}{S_{t+1}} \eta_t y_i x_i$  should be the update rule of  $W_t$ .

Decomposing  $w^{(t)}$  into  $S_t$  and  $W_t$  can significantly reduce the running time. Updating a number is much more timesaving than updating all parameters.