

## Problem Set 7, Nov 8, 2018 (Theory Questions Part)

### 2. Support Vector Machines using Coordinate Descent

1. The dual objective function that we have to optimize is the following :

$$\begin{aligned} & \underset{\alpha}{\text{maximize}} && f(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2\lambda} \alpha^\top Q \alpha \\ & \text{subject to} && \alpha \in [0, 1]^N \end{aligned}$$

where  $Q := \text{diag}(\mathbf{y}) \mathbf{X} \mathbf{X}^\top \text{diag}(\mathbf{y})$ . For computing coordinate update for one coordinate  $n$ , consider the following one variable sub-problem:

$$\begin{aligned} & \underset{\gamma \in \mathbb{R}}{\text{maximize}} && f(\alpha + \gamma \mathbf{e}_n) \\ & \text{subject to} && 0 \leq \alpha_n + \gamma \leq 1 \end{aligned}$$

where  $\mathbf{e}_n = [0, \dots, 1, \dots, 0]^\top$  (all zero vector except at the  $n^{\text{th}}$  position). Note that  $\gamma = 0$  (no need to update  $\alpha_n$ ) is an optimum iff  $\nabla_n^P f(\alpha) = 0$ , where  $\nabla_n$  denotes the  $n^{\text{th}}$  component of the gradient, and  $\nabla_n^P f(\alpha)$  is the projected gradient of  $f$ , projected onto the box or interval constraint.  $\nabla_n^P f(\alpha)$  can be computed as

$$\nabla_n^P f(\alpha) = \begin{cases} \nabla_n f(\alpha) & \text{if } 0 < \alpha_n < 1 \\ \min\{0, \nabla_n f(\alpha)\} & \text{if } \alpha_n = 0 \\ \max\{0, \nabla_n f(\alpha)\} & \text{if } \alpha_n = 1 \end{cases}$$

Note that,  $\nabla f(\alpha) = \mathbf{1} - \frac{1}{2\lambda} (Q + Q^\top) \alpha = \mathbf{1} - \frac{1}{\lambda} Q \alpha$  ( $Q$  is symmetric) and thus  $\nabla_n f(\alpha) = 1 - \frac{1}{\lambda} \mathbf{e}_n^\top Q \alpha$ . Simplifying  $f(\alpha + \gamma \mathbf{e}_n)$ , we get

$$\begin{aligned} f(\alpha + \gamma \mathbf{e}_n) &= (\alpha + \gamma \mathbf{e}_n)^\top \mathbf{1} - \frac{1}{2\lambda} (\alpha + \gamma \mathbf{e}_n)^\top Q (\alpha + \gamma \mathbf{e}_n) \\ &= \alpha^\top \mathbf{1} - \frac{1}{2\lambda} \alpha^\top Q \alpha + \gamma - \frac{1}{2\lambda} (\gamma^2 \mathbf{e}_n^\top Q \mathbf{e}_n + \gamma \alpha^\top Q \mathbf{e}_n + \gamma \mathbf{e}_n^\top Q \alpha) \\ &= f(\alpha) - \frac{\gamma^2}{2\lambda} Q_{nn} + \gamma(1 - \frac{1}{\lambda} \alpha^\top Q \mathbf{e}_n) \end{aligned}$$

Differentiating with respect to  $\gamma$  and equating to 0, we get :

$$\begin{aligned} & -\frac{\gamma}{\lambda} Q_{nn} + (1 - \frac{1}{\lambda} \alpha^\top Q \mathbf{e}_n) = 0 \\ & \gamma^* = \frac{\lambda}{Q_{nn}} (1 - \frac{1}{\lambda} \alpha^\top Q \mathbf{e}_n) \end{aligned}$$

Note that  $Q_{nn} = \mathbf{x}_n^\top \mathbf{x}_n y_n^2 = \mathbf{x}_n^\top \mathbf{x}_n$  and  $\alpha^\top Q \mathbf{e}_n = \sum_{i=1}^N \alpha_i Q_{i,n} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i^\top \mathbf{x}_n y_n$ . Using  $\mathbf{w}(\alpha) = \frac{1}{\lambda} \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$ , we get  $\alpha^\top Q \mathbf{e}_n = \lambda y_n \mathbf{w}^\top \mathbf{x}_n$  and thus

$$\gamma^* = \frac{\lambda}{\mathbf{x}_n^\top \mathbf{x}_n} (1 - y_n \mathbf{w}^\top \mathbf{x}_n)$$

We conclude

$$\begin{aligned} \alpha_n^{\text{new}} &= \alpha_n^{\text{old}} + \gamma^* \\ &= \alpha_n^{\text{old}} + \frac{\lambda}{\mathbf{x}_n^\top \mathbf{x}_n} (1 - y_n \mathbf{w}^\top \mathbf{x}_n) \end{aligned}$$

Since  $\alpha \in [0, 1]^N$ , we project  $\alpha_n$  as

$$\alpha_n^{\text{new}} := \min \left\{ \max \left\{ \alpha_n^{\text{old}} + \frac{\lambda}{\mathbf{x}_n^\top \mathbf{x}_n} (1 - y_n \mathbf{w}^\top \mathbf{x}_n), 0 \right\}, 1 \right\}$$

### 3. Kernels

1. Note that every term in the resulting polynomial is a product of kernels and that all these terms have positive coefficients. Hence the result follows by the two statements proved in class: namely that the positive sum of valid kernels is a valid kernel and that the product of valid kernels is a valid kernel.
2. We have  $\exp(x) = \sum_{i \geq 0} \frac{x^i}{i!}$ . We can hence apply the previous result concerning polynomials with positive coefficients and apply the limit.