

# ML HW-3

Q1 Monotonicity of loss + regularizer as the regularization parameter changes.

(A) By definition we have,

$$L(\tilde{\theta}^*) + \lambda r(\tilde{\theta}^*) \geq L(\theta^*) + \lambda r(\theta^*)$$

& similarly for  $\hat{\theta}$ ,

$$L(\theta^*) + \tilde{\lambda} r(\theta^*) \geq L(\hat{\theta}^*) + \tilde{\lambda} r(\hat{\theta}^*)$$

$\Rightarrow$  Subtracting the inequalities, we get the following results:

$$\Rightarrow \tilde{L}(\tilde{\theta}^*) + \lambda \mathcal{R}(\tilde{\theta}^*) - L(\tilde{\theta}^*) - \tilde{\lambda} \mathcal{R}(\tilde{\theta}^*)$$

$$\geq L(\theta^*) + \lambda \mathcal{R}(\theta^*) - L(\theta^*) - \lambda \mathcal{R}(\theta^*)$$

$$\Rightarrow (\lambda - \tilde{\lambda}) \mathcal{R}(\tilde{\theta}^*) \geq (\lambda - \tilde{\lambda}) \mathcal{R}(\theta^*),$$

f

moving  $\mathcal{R}(\theta^*)$  to the left hand side yields,

$$(\lambda - \tilde{\lambda})(\mathcal{R}(\tilde{\theta}^*) - \mathcal{R}(\theta^*)) \geq 0$$

But, since  $\tilde{\lambda} \geq \lambda$ , then

$$\mathcal{R}(\tilde{\theta}^*) - \mathcal{R}(\theta^*) \leq 0$$

Thus,

$$\boxed{\mathcal{R}(\tilde{\theta}^*) \leq \mathcal{R}(\theta^*)}$$

⑥ We can use the same trick by multiplying both sides of the inequalities by  $\frac{1}{\lambda}$ ,  $\frac{1}{\lambda}$  respectively then,

$$\left(\frac{1}{\lambda}\right) L(\tilde{\theta}^*) + \lambda(\tilde{\theta}^*) \geq \left(\frac{1}{\lambda}\right) L(\theta^*) + \lambda(\theta^*)$$

&

$$\left(\frac{1}{\lambda}\right) L(\theta^*) + \lambda(\theta^*) \geq \left(\frac{1}{\lambda}\right) L(\tilde{\theta}^*) + \lambda(\tilde{\theta}^*)$$

Subtracting both inequalities,

$$\Rightarrow \left(\frac{1}{\lambda} - \frac{1}{\lambda}\right) L(\tilde{\theta}^*) \geq \left(\frac{1}{\lambda} - \frac{1}{\lambda}\right) L(\theta^*)$$

$$\Rightarrow \left( \frac{1}{\lambda} - \frac{1}{\tilde{\lambda}} \right) \left( L(\tilde{\theta}^*) - L(\theta^*) \right) \geq 0$$

Since  $\lambda \leq \tilde{\lambda}$ , we have  $\left( \frac{1}{\lambda} - \frac{1}{\tilde{\lambda}} \right) \geq 0$ , so

$$L(\tilde{\theta}^*) - L(\theta^*) \geq 0$$

Thus,

$$L(\tilde{\theta}^*) \geq L(\theta^*)$$

Q2 MAP interpretation of regularized empirical loss minimization

A2 a)  $P(y|x_i; \theta) \sim N(\phi^T \theta, \sigma^2)$

$$P(\theta) \sim N(0, \sigma_0^2 I)$$

We deduce this:

$$P(y_i | x_i; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\phi^T \theta - y_i)^2}{2\sigma^2}}$$

$$P(\theta) = \frac{1}{\sqrt{2\pi\sigma_0^2 I}} e^{-\frac{\theta^2}{2\sigma_0^2 I}}$$

$\Rightarrow$  Use above two values in MAP formulation we get the following:

$$\text{minimize } \sum_{i=1}^n -\log \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\phi^T \theta - y_i)^2}{2\sigma^2}} - \log \frac{1}{\sqrt{2\pi\sigma_0^2 I}} e^{-\frac{\theta^2}{2\sigma_0^2 I}}$$

$\Rightarrow$  MAP depends on  $\theta$  we remove the constant terms and we get:

$$\text{minimize } \sum_{i=1}^n (\phi^T \theta - y_i)^2 + \frac{\sigma^2}{\sigma_0^2 I} \|\theta\|^2$$

$\Rightarrow$  This is of the form  $L(\theta) + \lambda R(\theta)$  where

$$\lambda = \frac{\sigma^2}{\sigma_0^2 I}$$

⑥

$$P(y|x; \theta) \sim N(\phi^T \theta, \sigma^2)$$

$$P(\theta) = \prod_{i=1}^m \frac{1}{\sigma} e^{-\frac{1}{\sigma^2} |\theta_i|^2}$$

$$-\log P(\theta) = -\lambda \log \left( \prod_{i=1}^m \frac{1}{\sigma} e^{-\frac{1}{\sigma^2} |\theta_i|^2} \right)$$

$$-\log P(\theta) = m \log \sigma + \sum_{j=1}^n \frac{|\theta_j|^2}{\sigma^2}$$

⇒ Taking from part ② & substituting in MAP formulation we get :

$$\Rightarrow \text{minimize } \sum_{i=1}^m (\phi^T \theta - y_i)^2 + \frac{\lambda \sigma^2}{\sigma^2} \sum_{j=1}^n |\theta_j|^2$$

⇒ This is of the form  $L(\theta) + \lambda R(\theta)$  where

$$\lambda = \frac{\lambda \sigma^2}{\sigma^2}$$

⑦  $P(y|x, \theta) = \Pr[y_a > 0]$  where

$$P(y|x, \theta) \sim N(\phi^T \theta, \sigma^2)$$

We can do the following:

$$P(y|x, \theta) = Pr[y \geq 0]$$

$$P(y|x, \theta) = \frac{1}{\sqrt{2\pi}} \int_0^{\infty} e^{-\frac{(u-y\phi^T\theta)^2}{2}} du$$

$$P(y|x, \theta) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{y\phi^T\theta} e^{-\frac{u^2}{2}} du$$

$$P(y|x, \theta) = \Phi(y\phi^T\theta)$$

Using part a we get

$$P(\theta) = \frac{1}{\sqrt{2\pi\sigma_0^2 I}} e^{-\frac{\theta^T\theta}{2\sigma_0^2 I}}$$

$\Rightarrow$  Substitution in MAP formula  
we get :

$$\text{minimize} -\log(\Phi(y\phi^T\theta)) + \frac{1}{2\sigma^2 I} \|\theta_j\|_2^2$$

$\Rightarrow$  This is of the form  $L(\theta) + \lambda R(\theta)$  where

$$\lambda = \frac{1}{2\sigma^2 I}$$

(Q)  $P(y_i, x_i | \theta) = e^{-y_i \phi^T \theta} / (1 + e^{-y_i \phi^T \theta})$

From (Q) the value of  $P(\theta)$  we can do  
following changes in our MAP formulation  
we get : minimize  $\sum_{i=1}^n y_i \phi^T \theta + \|\theta\|_2^2 / \lambda$

$\Rightarrow$  This is of the form  $L(\theta) + \lambda J(\theta)$  where

$$\lambda = \frac{1}{\alpha}$$

(Q4) Hand-Written Digits  
Classification

(A4) Given objective  $f$  :

$$\begin{aligned} & \text{minimize}_{\theta_1, \dots, \theta_K} \quad \frac{1}{n} \sum_{i=1}^n \max_c (\mathbf{x}_i^\top \theta_c - \mathbf{x}_i^\top \theta_{y_i} + 1_{y_i \neq c}) \\ & \quad + \lambda \sum_{j=1}^M \sqrt{\sum_{c=1}^K \theta_{jc}^2} \end{aligned} \quad \rightarrow ①$$

$\Rightarrow$  Randomly Select an  $\mathbf{x}_i$

$\Rightarrow$  Find subgradient at  $\theta_i^{(t)}$  by differentiating eq ①

$\Rightarrow \theta$  corresponding to correct class.

$$\nabla \theta_{y_i} L_i = - \left( \sum_{j \neq y_i} 1 (\theta_j^\top \mathbf{x}_i - \theta_{y_i}^\top \mathbf{x}_i + \alpha) \mathbf{x}_i \right)$$

$$\Rightarrow \theta_{y_i}^{(t+1)} = \theta_{y_i}^{(t)} - \gamma^{(t)} \mathbf{x}_i$$

$\Rightarrow \Theta$  corresponding to incorrect class

$$\nabla_{\Theta_j} L_i = \{ (\Theta_j^T x_i - \Theta_{y_i}^T x_i + \alpha) x_i$$

Alpha here is  $\underline{1}$ .

$$\Rightarrow \Theta_{y_i}^{(t+1)} = \Theta_{y_i}^{(t)} + \underline{\gamma^{(t)}} x_i$$

$\Rightarrow$  New  $\Theta$ (theta) is for proximal operation

$$\Theta^{t+1} = \Theta^t - \lambda \text{(Subgradient)}$$

$\Rightarrow$  New weight that is theta is given as:

$$\text{prox}(\Theta^t - \gamma^{(t)} \lambda \nabla L)$$

$\Rightarrow$  Subgradient is calculated by block

soft-thresholding on  $\Theta$ . The norm of each row of  $\Theta$  is calculated if it is less than  $\gamma^{(t)} \lambda$ , the entire row is set to 0, otherwise magnitude of the row decreases by the amount of  $\gamma^{(t)} \lambda$ .

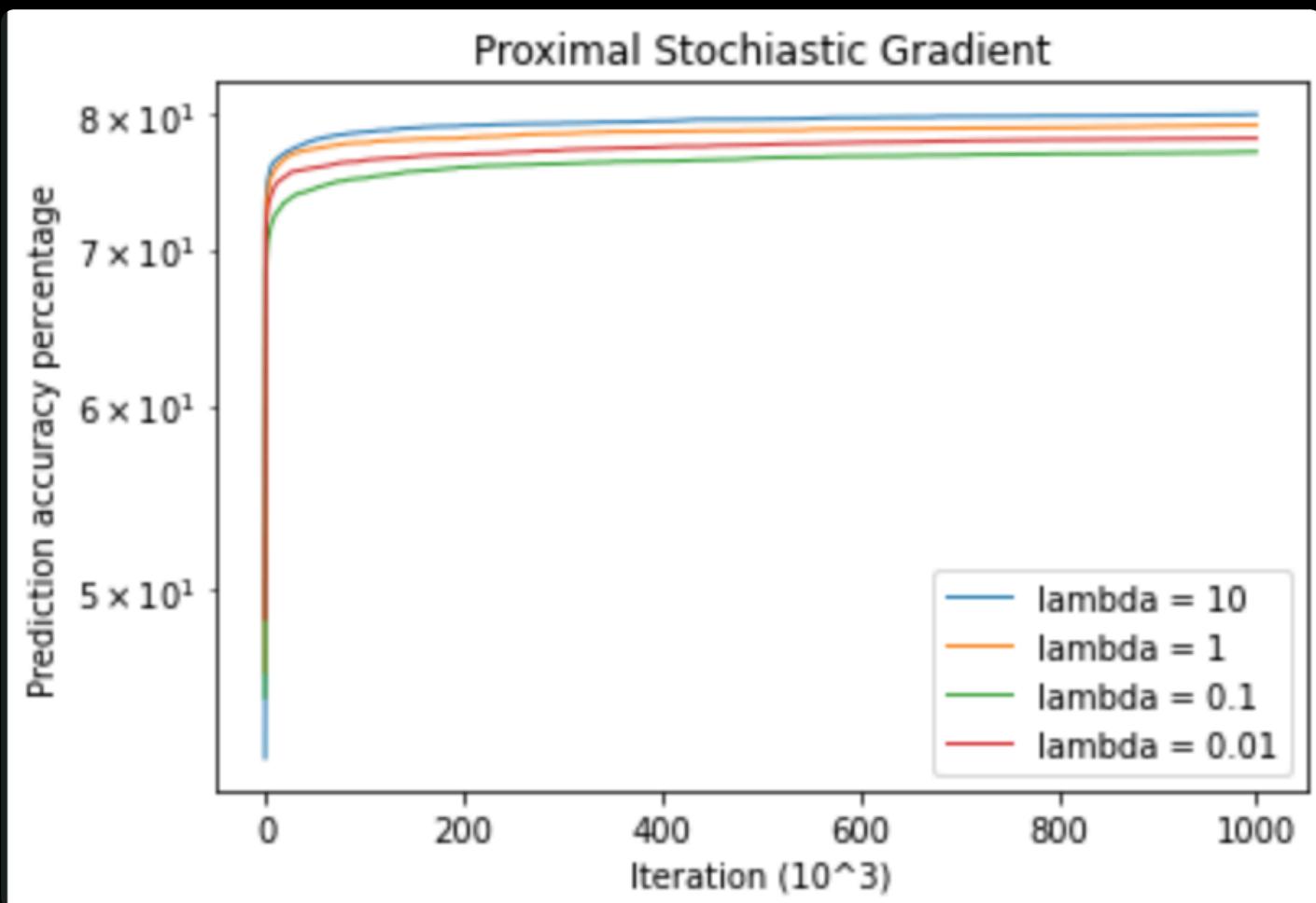
$$\Rightarrow \Theta_j^{(t+1)} = \begin{cases} 0 & : \text{if } \|v\|_2 \leq \gamma^{(t)} \lambda \\ \left(1 - \frac{\gamma^{(t)} \lambda}{\|v\|_2}\right) \Theta_j^{(t+1)} & : \text{otherwise} \end{cases}$$

$$\Rightarrow \text{where } V = \|\theta_j^{(t+1)}\|$$

⑥ Code attached in python file & run on kaggle

<https://www.kaggle.com/code/boltcoder/ml-hw3/notebook>

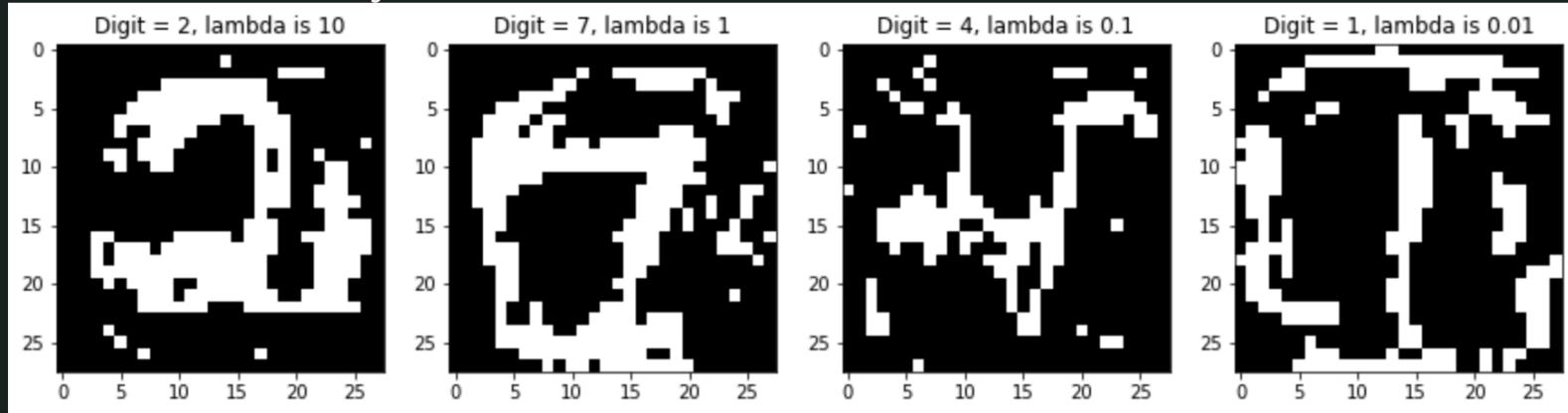
- c) The testing error - rate decreases as the algorithm progresses as shown in the figure for all  $\lambda$  values.  
→ Larger value of  $\lambda$  has higher accuracy value



d) For each  $\lambda$  we have a 0 solution & each 0 corresponds to 1 pixel.

Larger  $\lambda$  leads to fewer pixels being used for making predictions, where almost all pixels for  $\lambda=0.1$  &  $\lambda=0.01$  are used as shown in the figure.

Total No of features being discarded (Zeros) for lambda 10 : 257
Total No of features being discarded (Zeros) for lambda 1 : 170
Total No of features being discarded (Zeros) for lambda 0.1 : 114
Total No of features being discarded (Zeros) for lambda 0.01 : 93



(Q3) proximal Operator  
for the group lasso regularizer

A3 d)  $\|\theta_j\|_1 = \sqrt{\sum_{c=1}^k \theta_{j,c}^2}$

$\rightarrow$  If  $\theta_j \neq 0$ , then partial differentiation is given as follows:

$$\frac{\nabla \| \theta_j^* \|}{\nabla \theta_{j*}} = \frac{\theta_j^*}{\| \theta_j^* \|}$$

$$\Rightarrow \nabla \| \theta_j^* \| = \frac{\theta_j^*}{\| \theta_j^* \|}$$

$\rightarrow$  Else if  $\theta_j^* = 0$ ,

subgradient  $\nabla \| \theta_j^* \|$

is an element of  $\{z : \|z\| \leq 1\}$

$\rightarrow$  This holds because, by definition, in order of  $\nabla \| \theta_j^* \|$  to be a subgradient of ' $f$ ' we must have

$$\Rightarrow f(y) = \|y\| \geq f(x) + g^T(y-x) = g^T y$$

where  $f(x) = \|\theta_j\|$ ,

$$g = j \|\theta_j\| \quad \text{and} \quad x = \theta_j$$

$\Rightarrow \|y\| \geq g^T y$  holds true

only if the condition for  
 $g$  is  $\|g\| \leq 1$

$\therefore$  The Subdifferential  
 is given as follows:

For  $x \neq 0$ , unique subgradient

$$g = \frac{\theta_j}{\|\theta_j\|}$$

&

For  $x = 0$ , subgradient  $\theta_j$  is  
 any element of  $\{z : \|z\| \leq 1\}$

$$\textcircled{b} \quad \text{minimize } P \|\theta_j\| + \frac{1}{2} \|\theta_j - \tilde{\theta}_j\|^2$$

Differentiating equation w.r.t.  
 $\theta_j$  we get:

$$\Rightarrow P V + \theta_j - \tilde{\theta}_j \rightarrow \textcircled{1}$$

$$\text{where } V = \begin{cases} \frac{\theta_j}{\|\theta_j\|} & : \theta_j \neq 0 \\ z : \|z\| \leq 1 & : \theta_j = 0 \end{cases}$$

Now,

Equating \textcircled{1} with 0

we get:

$$P V + \theta_j - \tilde{\theta}_j = 0$$

$$\Rightarrow \theta_j = \tilde{\theta}_j - P V$$


---



---

