

CAP 6610 Homework 3

Siddhant Mittal (6061-8545)

1. (a) Given:

$$P(y/x; \theta) \sim N(\phi^T \theta, \sigma^2)$$
$$P(\theta) \sim N(0, \sigma_o^2 I)$$

From this we can deduce this:

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

Substituting the above two values in given explicit MAP formulation $\text{minimize} \sum_{i=1}^n -\log p(y_i/x_i, \theta) - \log p(\theta)$ we get the following:

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

Since MAP depends on θ we remove the constant terms and we get:

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

Which is of the form $L(\theta) + \lambda r(\theta)$ where

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

(b) Given:

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

$$P(\theta) = \prod_{j=1}^m \frac{1}{2a} \exp^{-\frac{|\theta_j|}{a}}$$
$$-\log p(\theta) = -\log \left(\prod_{j=1}^m \frac{1}{2a} \exp^{-\frac{|\theta_j|}{a}} \right)$$

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

Taking from part (a) $P(y_i/x_i; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\phi^T \theta - y_i)^2}{2\sigma^2}}$ and substituting in MAP formulation we get

$$\begin{aligned} \text{minimize} \quad & \sum_{i=1}^n (\phi^T \theta - y_i)^2 + \frac{2\sigma^2}{a} \sum_{j=1}^m |\theta_j| \\ \text{minimize} \quad & \sum_{i=1}^n (\phi^T \theta - y_i)^2 + \frac{2\sigma^2}{a} \|\theta_j\|_1 \end{aligned}$$

Which is of the form $L(\theta) + \lambda r(\theta)$ where

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

(c) Given:

$$p(y/x, \theta) = Pr[yu \geq 0] \text{ where } p(u/x, \theta) \sim N(\phi^T \theta, \sigma^2)$$

We can do the following:

$$p(y/x, \theta) = Pr[yu \geq 0]$$

$$p(y/x, \theta) = \frac{1}{\sqrt{2\pi}} \int_0^\infty \exp\left(-\frac{(u - y\phi^T \theta)^2}{2}\right) du$$

$$p(y/x, \theta) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{y\phi^T \theta} \exp\left(-\frac{u^2}{2}\right) du$$

$$p(y/x, \theta) = \Phi(y\phi^T \theta)$$

Using part a we get

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

Substituting in MAP formula we get

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

Which is of the form $L(\theta) + \lambda r(\theta)$ where

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

(d) Given:

$$p(y_i, x_i; \theta) = \frac{e^{-y_i \phi^T \theta}}{1 + e^{-y_i \phi^T \theta}}$$

and taking from part (b) the value of $p(\theta)$ we can do following changes in our MAP formulation and ignoring constant terms we get:

$$\text{minimize} \sum_{i=1}^n y_i \phi^T \theta + \frac{\|\theta\|_1}{a}$$

Which is of the form $L(\theta) + \lambda r(\theta)$ where

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

2. (a) We know that

$$\text{Prox}_f(\theta) = \text{argmin} f(\theta) + \frac{1}{2} \|\theta - \theta_1\|^2$$

To find proximal mapping for θ we need to minimize the above equation that means we can write

$$\theta = \text{Prox}_f(\theta_1)$$

Differentiating RHS and using sum rule of sub differentiation we get

$$\frac{\partial}{\partial t} [f(\theta) + \frac{1}{2} \|\theta - \theta_1\|^2] \in 0$$

$$\partial f(\theta) + \frac{1}{2} 2(\theta - \theta_1)(1 - 0) \in 0$$

$$\theta_1 - \theta \in \partial f(\theta)$$

Since $\theta = \text{Prox}_f(\theta_1)$ we put that in

$$\theta_1 - \text{Prox}_f(\theta_1) \in \partial(\text{Prox}_f(\theta_1))$$

(b) Given $g_1 \in \partial f(\theta_1)$ and $g_2 \in \partial f(\theta_2)$ we have to prove $(g_1 - g_2)^T(\theta_1 - \theta_2) \geq 0$. By definition of subgradient of $f(\theta_1)$ and $f(\theta_2)$ we have following

$$\begin{aligned} f(\theta) &\geq f(\theta_1) + g_1^T(\theta - \theta_1) \forall \theta \\ f(\theta) &\geq f(\theta_2) + g_2^T(\theta - \theta_2) \forall \theta \end{aligned}$$

Since it is for all θ this can also be written as

$$\begin{aligned} f(\theta_2) &\geq f(\theta_1) + g_1^T(\theta_2 - \theta_1) \forall \theta \\ f(\theta_1) &\geq f(\theta_2) + g_2^T(\theta_1 - \theta_2) \forall \theta \end{aligned}$$

Adding above two equations we get

$$g_1^T(\theta_2 - \theta_1) + g_2^T(\theta_1 - \theta_2) \leq 0$$

Simplifying gives

$$(g_1 - g_2)^T(\theta_1 - \theta_2) \geq 0$$

(c)

From part (a) we have $\theta_1 - \text{Prox}_f(\theta_1) \in \partial(\text{Prox}_f(\theta_1))$ and from part (b) we have $(g_1 - g_2)^T(\theta_1 - \theta_2) \geq 0$. Putting values of θ_1 and θ_2 in form of part (a) in part (b) equation we get

Putting

$$\begin{aligned} g_1 &= \theta_1 - \text{Prox}_f(\theta_1) \\ g_2 &= \theta_2 - \text{Prox}_f(\theta_2) \end{aligned}$$

$$(\theta_1 - \text{Prox}_f(\theta_1) - (\theta_2 - \text{Prox}_f(\theta_2)))^T(\theta_1 - \theta_2) \geq 0$$

By simplifying this we get

$$(\text{Prox}_f(\theta_1) - \text{Prox}_f(\theta_2))^T(\theta_1 - \theta_2) \geq \|\text{Prox}_f(\theta_1) - \text{Prox}_f(\theta_2)\|^2$$

(d)

By applying Cauchy-Schwartz inequality property $a^T b \geq \|a\| \|b\|$ is $b \geq \|a\|$

We can modify the part equation to get nonexpansiveness property.

$$\|\theta_1 - \theta_2\| \geq \|\text{Prox}_f(\theta_1) - \text{Prox}_f(\theta_2)\|$$

3. (a) Following algorithm is derived for solving this problem.

Given objective function.

$$\text{minimize}_{\theta_1, \dots, \theta_k} \frac{1}{n} \sum_{i=1}^n \max_c (x_i^T \theta_c - x_i^T \theta_{y_i} + 1_{y_i \neq c}) \quad (1)$$

$$+ \lambda \sum_{j=1}^m \int \sum_{c=1}^K \theta_{jc}^2$$

→ Randomly choose select a x_i .

→ Find the ^{sub}gradient at x_i by differentiation equation (1)

→ Find θ corresponding to correct class.

$$\nabla_{\theta_{y_i}} L_i = - \left(\sum_{y \neq y_i} 1 (\theta_j^T x_i - \theta_{y_i}^T x_i + \text{delta}) \right) x_i$$

θ corresponding to incorrect class.

$$\nabla_{\theta_j} L_i = 1 (\theta_j^T x_i - \theta_{y_i}^T x_i + \text{delta}) x_i$$

delta here is one

→ new θ (theta) is for proximal operation
 $\theta^{t+1} \leftarrow \theta^t - \lambda (\text{subgradient})$

→ New weight that is theta is

$$\text{prox}(\theta^t - y^{(t)} \lambda \nabla L)$$



Scanned with
CamScanner

→ These steps are repeated 106 iterations for 4 lambda values $\{10, 1, 0.1, 0.01\}$

→ $y(t)$ is the step size which is diminishing step size at each iteration, step size is $\frac{1}{t+1}$ (1st iteration is the start)

Proximal operator calculation:

Prox of $L_{2,1}$ Mixed Norm

The problem is given by:

$$\arg \min_X \frac{1}{2} \|X - Y\|_F^2 + \lambda \|X\|_{2,1}$$

Where $X, Y \in \mathbb{R}^{m \times n}$.

Again, this can be decomposed into working on each column of X separately:

$$\begin{aligned} \arg \min_X \frac{1}{2} \|X - Y\|_F^2 + \lambda \|X\|_{2,1} &= \arg \min_X \sum_i \frac{1}{2} \|X_{:,i} - Y_{:,i}\|_2^2 + \sum_i \lambda \|X_{:,i}\|_2 \\ &= \arg \min_X \left(\frac{1}{2} \|X_{:,1} - Y_{:,1}\|_2^2 + \lambda \|X_{:,1}\|_2^2 \right) \\ &\quad + \left(\frac{1}{2} \|X_{:,2} - Y_{:,2}\|_2^2 + \lambda \|X_{:,2}\|_2^2 \right) \\ &\quad + \dots \\ &\quad + \left(\frac{1}{2} \|X_{:,n} - Y_{:,n}\|_2^2 + \lambda \|X_{:,n}\|_2^2 \right) \end{aligned}$$

Each term in the brackets is independent [Prox Function of \$L_2\$ Norm](#).

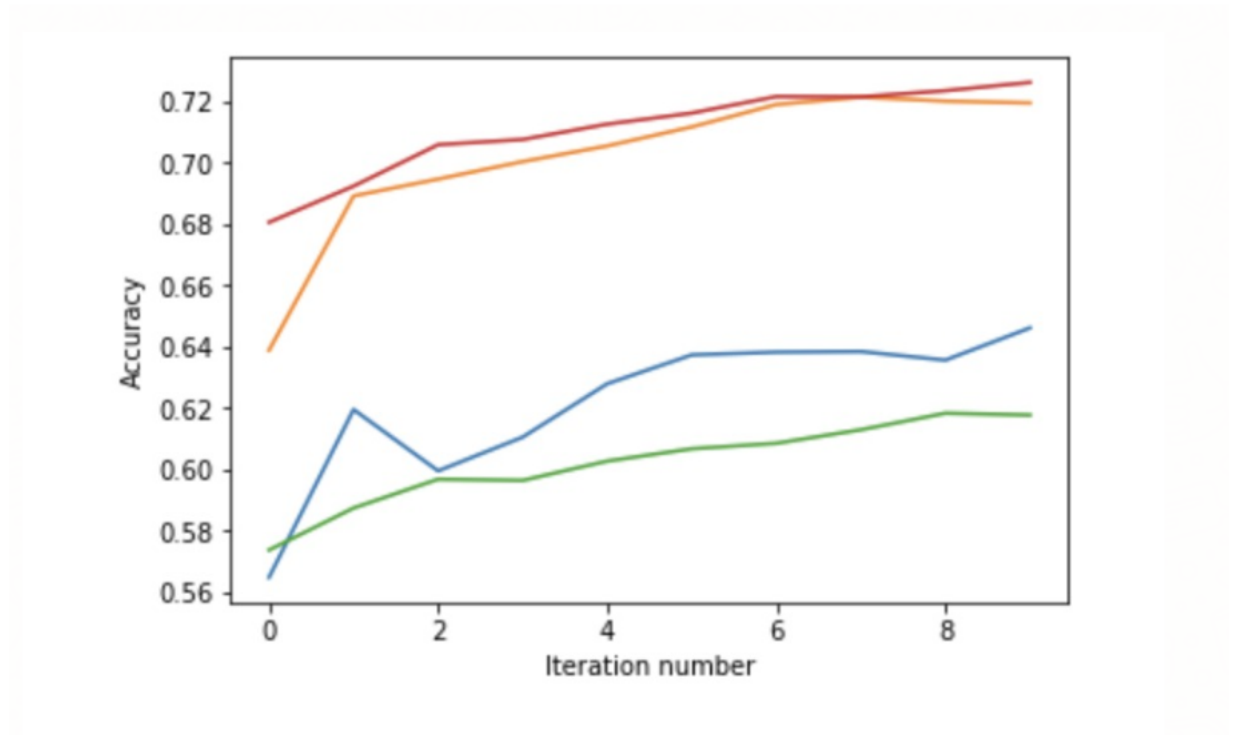
Hence the solution is given by:

$$\hat{X} = \arg \min_X \frac{1}{2} \|X - Y\|_F^2 + \lambda \|X\|_{2,1}$$

$$\text{Where } \hat{X}_{:,i} = Y_{:,i} \left(1 - \frac{\lambda}{\max(\|Y_{:,i}\|_2, \lambda)} \right)$$

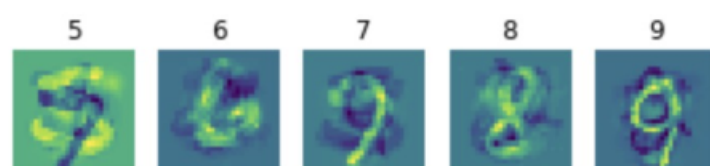
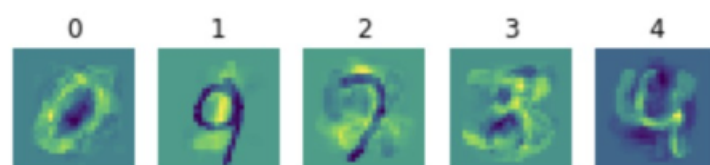
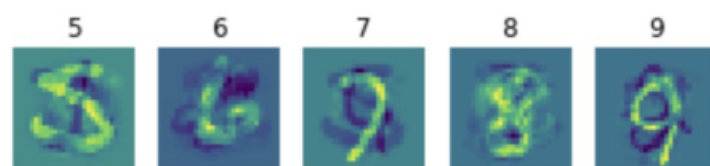
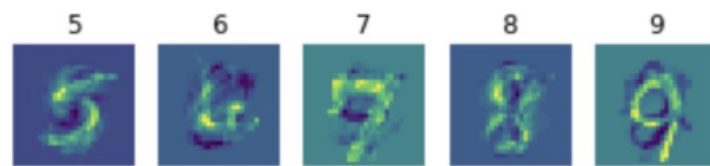
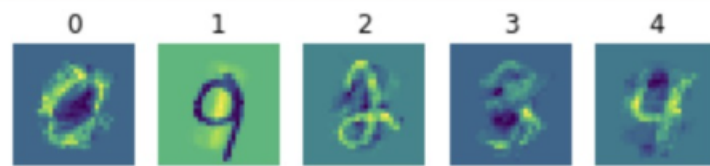
(c) Red line is $\lambda = 0.01$, Green line is $\lambda = 0.1$, Orange line is

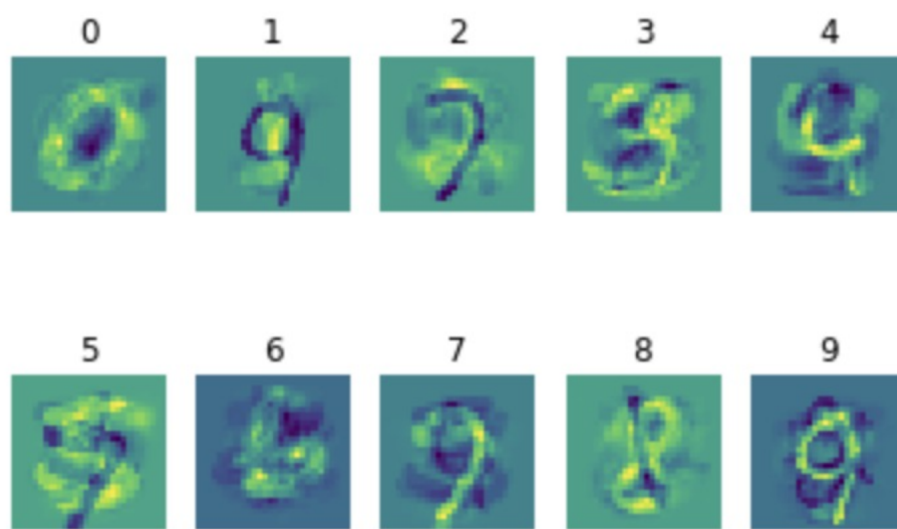
$\lambda = 1$, Blue line is $\lambda = 10$



(d) Yes it true that a large λ leads to a more sparse solution. Below shows exactly the same.

Weights are the order of λ 10,1,0,1,0.01 respectively





0 1 2 3 4 5 6 7 8 9