

$$1.a) \min_{\theta} \sum_{i=1}^n -\log p(y_i | x_i, \theta) - \log p(\theta)$$

$$p(y|x, \theta) \sim \mathcal{N}(\phi^T \theta, \sigma^2)$$

$$p(\theta) \sim \mathcal{N}(0, \sigma_0^2 I)$$

Applying distribution to the function,

$$\min_{\theta} \sum_{i=1}^n -\log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \phi^T \theta)^2}{2\sigma^2}\right) \right] - \log \left[\frac{1}{\sqrt{|2\pi\sigma_0^2 I|}} \exp\left(-\frac{1}{2}(\theta-0)^T (\sigma_0^2 I)^{-1} (\theta-0)\right) \right]$$

\Downarrow

$$\frac{1}{(2\pi\sigma_0^2)^{d/2}} \exp\left(-\frac{1}{2\sigma_0^2} \|\theta\|^2\right)$$

$$\min_{\theta} \sum_{i=1}^n \left(\underbrace{\frac{1}{2} \log(2\pi\sigma^2)}_{\text{const}} + \frac{(y_i - \phi^T \theta)^2}{2\sigma^2} + \underbrace{\frac{d}{2} \log(2\pi\sigma^2)}_{\text{const}} + \frac{1}{2\sigma_0^2} \|\theta\|^2 \right)$$

$$\min_{\theta} \sum_{i=1}^n \frac{(y_i - \phi^T \theta)^2}{2\sigma^2} + \frac{n}{2\sigma_0^2} \|\theta\|^2 = \min_{\theta} \underbrace{\sum_{i=1}^n \frac{(y_i - \phi^T \theta)^2}{2\sigma^2}}_{\text{quadratic loss}} + \underbrace{\frac{n}{2\sigma_0^2} \|\theta\|^2}_{\text{ridge}}$$

$$\text{So, } \boxed{\lambda = \frac{n\sigma^2}{\sigma_0^2}}$$

$$b) p(y|x, \theta) \sim \mathcal{N}(\phi^T \theta, \sigma^2) \quad p(\theta) = \prod_{j=1}^m \frac{1}{2a} \exp\left(-\frac{|0_j|}{a}\right)$$

$$\min_{\theta} \sum_{i=1}^n \left(\frac{1}{2} \log(2\pi\sigma^2) + \frac{(y_i - \phi^T \theta)^2}{2\sigma^2} + \underbrace{\sum_{j=1}^m \left(-\log(2a) - \frac{|0_j|}{a} \right)}_{\text{const}} \right)$$

$$\min_{\theta} \sum_{i=1}^n \left(\frac{(y_i - \phi^T \theta)^2}{2\sigma^2} \right) + \frac{n}{a} \|\theta\|_1 \quad \left[\sum_{j=1}^m |0_j| \rightarrow \ell_1 \text{ norm} \right]$$

$$\text{So, } \min_{\theta} \underbrace{\sum_{i=1}^n \frac{(y_i - \phi^T \theta)^2}{2\sigma^2}}_{\text{quadratic loss}} + \underbrace{\frac{2n\sigma^2}{a} \|\theta\|_1}_{\text{lasso}}$$

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

c) $p(y|x, \theta) = P_c[yu \geq 0]$ where $y = \pm 1$, $p(u|x, \theta) \sim \mathcal{N}(\phi^T \theta, \sigma^2)$
 $p(\theta) \sim \mathcal{N}(0, \sigma_0^2 \mathbf{I})$

$$p(y|x, \theta) = P_c[yu \geq 0] = \frac{1}{\sqrt{2\pi}\sigma^2} \int_0^\infty \exp\left(-\frac{1}{2} \frac{(u - \phi^T \theta)^2}{\sigma^2}\right) du$$

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

So, minimizing function $\min_{\theta} \sum_{i=1}^n \underbrace{-\log \Phi(y_i \phi_i^T \theta)}_{\text{probit loss}} - \underbrace{\frac{\eta}{2\sigma_0^2} \|\theta\|^2}_{\text{from derivation in a.}}$

$\lambda = \frac{\eta}{2\sigma_0^2}$

d) $p(y|x, \theta) = \frac{1}{1 + \exp(-y\phi^T \theta)}$ where $y = \pm 1$, $p(\theta) = \prod_{j=1}^m \frac{1}{2a} \exp\left(-\frac{|a_j|}{a}\right)$

So, minimizing function, $\min_{\theta} \sum_{i=1}^n \left(-\log\left(\frac{1}{1 + \exp(-y_i \phi_i^T \theta)}\right) - \underbrace{\sum_{j=1}^m \left(-\log a - \frac{|a_j|}{a}\right)}_{\text{const}} \right)$

$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i \phi_i^T \theta)) + \sum_{i=1}^n \|\theta\|_1$$

$$\Rightarrow \min_{\theta} \underbrace{\sum_{i=1}^n \log(1 + \exp(-y_i \phi_i^T \theta))}_{\text{logistic loss}} + \underbrace{\sum_{i=1}^n \|\theta\|_1}_{\lambda}$$

$$2. \quad \text{prox}_f(\theta_1) = \arg \min_{\theta} f(\theta) + \frac{1}{2} \|\theta - \theta_1\|^2$$

By definition, $\text{prox}_f(\theta_1)$ minimizes the problem. From optimality,

$$0 \in \partial f(\text{prox}_f(\theta_1)) + \text{prox}_f(\theta_1) - \theta_1$$

a)

$$\Rightarrow \theta_1 - \text{prox}_f(\theta_1) \in \partial f(\text{prox}_f(\theta_1)) \text{ — Proved. } \textcircled{1}$$

b)

$$g_1 \in \partial f(\theta_1) \quad g_2 \in \partial f(\theta_2)$$

g_1 is a subgradient of f at θ_1 , g_2 is a subgradient of f at θ_2 .

Now, subgradient is a global underestimator.

$$f(\theta_2) \geq f(\theta_1) + g_1^T(\theta_2 - \theta_1) \text{ — (i)}$$

$$f(\theta_1) \geq f(\theta_2) + g_2^T(\theta_1 - \theta_2) \text{ — (ii)}$$

adding (i) and (ii), $0 \geq g_1^T(\theta_2 - \theta_1) + g_2^T(\theta_1 - \theta_2)$

$$\text{So, } (g_1 - g_2)^T(\theta_1 - \theta_2) \geq 0 \text{ — Proved. } \textcircled{11}$$

c) From a and b,

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

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

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

$$d) \quad \|\text{prox}_f(\theta_1) - \text{prox}_f(\theta_2)\|^2 \leq \|\text{prox}_f(\theta_1) - \text{prox}_f(\theta_2)\| \cdot \|\theta_1 - \theta_2\|$$

(Applying Cauchy Schwarz)

$$\therefore \|\text{prox}_f(\theta_1) - \text{prox}_f(\theta_2)\| \leq \|\theta_1 - \theta_2\|$$

3. a)

$$\underset{\theta_1, \dots, \theta_k}{\text{minimize}} \quad \sum_{i=1}^m \max_c (\alpha_i^T \theta_c - \alpha_i^T \theta_{y_i} + 1) + \lambda \sum_{j=1}^n \sqrt{\sum_{c=1}^k \theta_{jc}^2}$$

At iteration t ,

For each class $c = 1 \dots k$

① Pick a sample i from class c train data

② Pick gradient based on data sample i ,

$$\text{gradient based on } \nabla \theta_{y_i} L_i = - \left(\sum_{c \neq y_i} 1(\alpha_i^T \theta_c - \alpha_i^T \theta_{y_i} + 1 > 0) \right) \alpha_i$$

$$\nabla \theta_j L_i = 1(\alpha_i^T \theta_c - \alpha_i^T \theta_{y_i} + 1 > 0) \alpha_i$$

Update θ accordingly.

$$[\delta^{(t)} = \frac{1}{t}] \quad \theta_{y_i} = \theta - \frac{1}{t} \nabla \theta_{y_i} L_i$$

$$\theta_j = \theta_j - \frac{1}{t} \nabla \theta_j L_i$$

③ Apply Proximal operator,

For $j = 1$ to m ,

$$\theta_j^{(t+1)} = \text{Prox}_{\delta \lambda} \theta_j^{(t)}, \text{ where } \text{Prox}_{\delta \lambda}(\theta) = \begin{cases} (1 - \frac{\delta \lambda}{\|\theta\|}) \theta & \text{if } \|\theta\| > \delta \lambda \\ 0 & \text{if } \|\theta\| \leq \delta \lambda \end{cases} \quad \text{--- (iv)}$$

$$\text{① } \|\theta_j\| > \delta \lambda$$

$$\text{② } \|\theta_j\| \leq \delta \lambda$$

d) Upon the implementation of our algo, the feature vector 784, we see we are getting different count of 0 vectors depending on regularization parameter.

① $\lambda = 10$, 249 features are discarded.

② $\lambda = 1$, 179 features are discarded.

③ $\lambda = 0.1$, 116 features are discarded.

④ $\lambda = 0.01$, 91 features are discarded.