

problem 2.1)

$$a) f(\omega_2) - f(\omega_1) = \frac{1}{N} \sum_i f_i(\omega_2) - f_i(\omega_1) + \lambda (\|\omega_2\|^2 - \|\omega_1\|^2)$$

$$\leq \frac{1}{N} \sum_i f_i(\omega_2) - f_i(\omega_1) + \lambda \|\omega_2 - \omega_1\|^2$$

$$= \frac{1}{N} \sum_i \log \frac{1 + e^{-y_i \omega_2^T x_i}}{1 + e^{-y_i \omega_1^T x_i}} + \lambda \|\omega_2 - \omega_1\|^2$$

convexity of  
 $\log(\cdot)$

$$\leq \log \frac{1}{N} \sum_i \frac{1 + e^{-y_i \omega_2^T x_i}}{1 + e^{-y_i \omega_1^T x_i}} + \lambda \|\omega_2 - \omega_1\|^2$$

$$= \log \frac{1}{N} \sum_i \frac{e^{y_i \omega_1^T x_i} + e^{-y_i (\omega_2 - \omega_1)^T x_i}}{1 + e^{y_i \omega_1^T x_i}} + \lambda \|\omega_2 - \omega_1\|^2$$

$$= \log \frac{1}{N} \sum_i \frac{1 + e^{-y_i (\omega_2 - \omega_1)^T x_i}}{1 + e^{y_i \omega_1^T x_i}} + \lambda \|\omega_2 - \omega_1\|^2$$

$$\leq \log \frac{1}{N} \sum_i e^{-y_i (\omega_2 - \omega_1)^T x_i} + \lambda \|\omega_2 - \omega_1\|^2 = A$$

Now from Cauchy-Schwarz we have:

$$|y_i (\omega_2 - \omega_1)^T x_i| \leq \underbrace{|y_i|}_{=1} \underbrace{\|\omega_2 - \omega_1\|_2}_{\leq c} \underbrace{\|x_i\|_2}_{\leq c}$$

$$\leq c \|\omega_2 - \omega_1\|^2$$

$$\text{so } A \leq \log \frac{1}{N} \sum_i e^{c \|\omega_2 - \omega_1\|^2} + \lambda \|\omega_2 - \omega_1\|^2$$

$$= (\lambda + c) \|\omega_2 - \omega_1\|^2$$

similar steps hold for  $f(\omega_1) - f(\omega_2)$  and

$$f(\omega_1) - f(\omega_2) \leq (\lambda + c) \|\omega_2 - \omega_1\|^2$$

so  $|f(\omega_1) - f(\omega_2)| \leq (\lambda + c) \|\omega_2 - \omega_1\|^2$ . so  $f(\cdot)$  is  $(\lambda + c)$ -Lipschitz

$$\text{and } B = \lambda + c$$

2.1)

$$b) \quad f_i = \log 1 + e^{-y_i \omega^T x_i} \quad \nabla f_i = \left( - \frac{y_i}{e^{y_i \omega^T x_i} + 1} \right) x_i$$

Since  $f_i$  is twice differentiable, a smoothness is equivalent to  $\nabla^2 f(\omega) \leq L I$  for some  $L > 0$ .

$$\nabla^2 f_i = x_i x_i^T h''(y_i \omega^T x_i) \quad \text{where } h(z) = \log(1 + e^{-z}) \quad \text{and } h'(z) = \frac{-e^{-z}}{1 + e^{-z}} \quad \text{and } h''(z) = \frac{e^{-z}}{(1 + e^{-z})^2}$$

$$\text{Now } h''(z) \leq 1 \quad \text{so } \delta_{\max}(\nabla^2 f_i) \leq \delta_{\max}(x_i x_i^T)$$

$$\text{So } \nabla^2 f(\omega) \leq \underbrace{\delta_{\max}(x_i x_i^T)}_{\text{A small suggested } L} I \quad \text{so it is } L\text{-smooth.}$$

Similarly  $f(\cdot)$  is also smooth.

$$\nabla f = \frac{1}{N} \sum_i \nabla f_i + 2\lambda \omega \quad \nabla^2 f = \frac{1}{N} \sum_i \nabla^2 f_i + 2\lambda I$$

$$\delta_{\max}(\nabla^2 f) \leq \underbrace{\delta_{\max}\left(\frac{1}{N} \sum_i x_i x_i^T\right) + 2\lambda}_{\triangleq \delta^*} \Rightarrow \nabla^2 f(\omega) \leq \delta^* I \quad \text{so } f(\omega) \text{ is } \delta^*\text{-smooth.}$$

c) Being  $\mu$ -strongly convex is equivalent to  $\nabla^2 f(\omega) \geq \mu I$  for twice differentiable  $f(\cdot)$ .

$$\nabla^2 f = \frac{1}{N} \sum_i \nabla^2 f_i + 2\lambda I \quad \text{note that } \nabla^2 f_i = x_i x_i^T h''(y_i \omega^T x_i) \quad \text{and } h''(z) \geq 0$$

$$\begin{aligned} \text{so } \min_{\substack{u \\ \|u\|^2=1}} u^T \nabla^2 f u &= \min_{\|u\|^2=1} \frac{1}{N} \sum_i u^T x_i x_i^T u h''(y_i \omega^T x_i) + 2\lambda \|u\|^2 \\ &= \min_{\|u\|^2=1} \frac{1}{N} \sum_i \underbrace{(u^T x_i)^2}_{\geq 0} \underbrace{h''(y_i \omega^T x_i)}_{\geq 0} + 2\lambda \end{aligned}$$

$$\geq 2\lambda$$

$$\text{so } \nabla^2 f(\omega) \geq 2\lambda I \quad \rightarrow \text{a large } \mu$$

so  $f(\omega)$  is  $2\lambda$ -strongly convex

$$2.2) \quad (1) \quad \nabla f(\omega_k)^T E_{\xi_k} [g(\omega_k; \xi_k)] \geq c \|\nabla f(\omega_k)\|_2^2 \quad \exists c_0 \geq c > 0 \quad \forall k$$

$$(2) \quad \|E_{\xi_k} [g(\omega_k; \xi_k)]\|_2 \leq c_0 \|\nabla f(\omega_k)\|_2$$

$$(3) \quad \text{Var}_{\xi_k} [g(\omega_k; \xi_k)] \leq M + M_V \|\nabla f(\omega_k)\|_2^2 \quad \exists M, M_V \geq 0 \quad \forall k$$

$$E_{\xi_k} [\|g(\omega_k; \xi_k)\|^2] = \text{Var}_{\xi_k} [g(\omega_k; \xi_k)] + \|E[g(\omega_k; \xi_k)]\|^2$$

$$\leq M + M_V \|\nabla f(\omega_k)\|_2^2 + c_0^2 \|\nabla f(\omega_k)\|_2^2$$

$$\text{so } \alpha = M \quad \beta = M_V + c_0^2$$



## Homework Assignment #2

### Problem 2.3

$$\begin{aligned} * E E_K [F(w_{K+1}) - F(w)] &\leq -\mu \alpha_K \|\nabla F(w_K)\|_2^2 + \frac{1}{2} \alpha_K^2 L E E_K [\|g(w_K, \varepsilon_K)\|_2^2] \quad (1) \\ &\leq -(\mu - \frac{1}{2} \alpha_K L M_G) \alpha_K \|\nabla F(w_K)\|_2^2 + \frac{1}{2} \alpha_K^2 K L M \quad (2) \end{aligned}$$

since  $\alpha_K L M_G \leq \mu$  for all  $K \in \mathbb{N}$  by assumption and  $E E_K [\|g(w_K, \varepsilon_K)\|_2^2] \leq \alpha + \beta \|\nabla f(w_K)\|_2^2$  from problem (2.2)

\* Step size requirements:  $\sum_{K=1}^{\infty} \alpha_K = \infty \quad (a)$   
 $\sum_{K=1}^{\infty} \alpha_K^2 < \infty \quad (b) \quad (3)$

\* summing both sides of the inequality in (2) for  $K \in \{1, \dots, K\}$  gives

$$F(w) - E[F(w_1)] \leq E[F(w_{K+1})] - E[F(w_1)] \quad (4)$$

$$\leq -\frac{1}{2} \mu \sum_{K=1}^K \alpha_K E[\|\nabla F(w_K)\|_2^2] + \frac{1}{2} L M \sum_{K=1}^K \alpha_K^2 \quad (5)$$

\* dividing by  $\mu/2$  and rearranging the terms, we obtain

$$\sum_{K=1}^K \alpha_K E[\|\nabla F(w_K)\|_2^2] \leq \frac{2(E[F(w_1)] - F(w))}{\mu} + \frac{L M}{\mu} \sum_{K=1}^{\infty} \alpha_K^2 \quad (7)$$

\* let  $A_K := \sum_{K=1}^K \alpha_K$

\* 3(b) implies that the right-hand side of this inequality converges to a finite limit when  $K$  increases proving that:

$$\lim_{K \rightarrow \infty} E \left[ \sum_{K=1}^{\infty} \alpha_K \|\nabla F(w_K)\|_2^2 \right] < \infty$$

\* 3(a) ensures that  $A_K \rightarrow \infty$  as  $K \rightarrow \infty$ , proving that

$$E \left[ \frac{1}{A_K} \sum_{K=1}^{\infty} \alpha_K \|\nabla F(w_K)\|_2^2 \right] \xrightarrow{K \rightarrow \infty} 0$$

### References

L. Bottou, F.E. Curtis and J. Nocedal, "Optimization methods for large-scale machine learning," SIAM Review, 2018.