

# Homework #2

( Due: Apr 19 )

## 1 Theory

### 1.1 Gaussian Distributions

To prove :

$$E[x] = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} x dx \quad (1)$$

Substituting

$$t = \frac{(x - \mu)}{\sqrt{2}\sigma}$$

we get

$$E[x] = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} (\sqrt{2}\sigma t + \mu)(\sqrt{2}\sigma) e^{-t^2} dt \quad (2)$$

$$\implies E[x] = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} (\sqrt{2}\sigma t + \mu) e^{-t^2} dt \quad (3)$$

$$\implies E[x] = \frac{1}{\sqrt{\pi}} \left( \sqrt{2}\sigma \int_{-\infty}^{\infty} t e^{-t^2} dt + \mu \int_{-\infty}^{\infty} e^{-t^2} dt \right) \quad (4)$$

$$\implies E[x] = \frac{1}{\sqrt{\pi}} \left( \sqrt{2}\sigma \left[ \frac{-e^{-t^2}}{2} \right]_{-\infty}^{\infty} + \mu \int_{-\infty}^{\infty} e^{-t^2} dt \right) \quad (5)$$

$$\implies E[x] = \frac{\mu}{\sqrt{\pi}} \left( \int_{-\infty}^{\infty} e^{-t^2} dt \right) \quad (6)$$

Directly using the result of  $\left( \int_{-\infty}^{\infty} e^{-t^2} dt \right) = \sqrt{\pi}$ , since the function is non-integrable in one dimension.

$$\implies E[x] = \frac{\mu}{\sqrt{\pi}} (\sqrt{\pi}) = \mu \quad (7)$$

Since,

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = 1 \quad (8)$$

Let  $t = \sigma^2$ , differentiating the above with respect to  $t$  we get,

$$\Rightarrow \int_{-\infty}^{\infty} \frac{-1}{2\sqrt{2\pi t^{\frac{3}{2}}}} e^{\frac{-(x-\mu)^2}{2t}} dx + \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi t}} e^{\frac{-(x-\mu)^2}{2t}} \left( \frac{(x-\mu)^2}{2t^2} \right) dx = 0 \quad (9)$$

$$\Rightarrow \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi t^{\frac{3}{2}}}} e^{\frac{-(x-\mu)^2}{2t}} dx = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi t}} e^{\frac{-(x-\mu)^2}{2t}} \left( \frac{(x-\mu)^2}{t^2} \right) dx \quad (10)$$

$$\Rightarrow \frac{1}{t} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi t}} e^{\frac{-(x-\mu)^2}{2t}} dx = \frac{1}{t^2} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi t}} e^{\frac{-(x-\mu)^2}{2t}} (x^2 - 2\mu x + \mu^2) dx \quad (11)$$

$$\Rightarrow t = E[x^2] - 2\mu E[x] + \mu^2 \quad (12)$$

$$\begin{aligned} \Rightarrow E[x^2] &= t + 2\mu E[x] - \mu^2 \\ &= t + 2\mu^2 - \mu^2 \\ &= \sigma^2 + \mu^2 \end{aligned} \quad (13)$$

Therefore  $\text{var}(x) = E[x^2] - (E[x])^2 = \sigma^2 + \mu^2 - \mu^2 = \sigma^2$ .

## 1.2 Strongly convex function

For a strongly convex function with parameter  $\lambda$  we can write

$$f(\alpha u + (1-\alpha)w) \leq \alpha f(u) + (1-\alpha)f(w) - \frac{\lambda}{2}\alpha(1-\alpha)\|u-w\|^2 \quad (14)$$

$$\Rightarrow f(\alpha u + (1-\alpha)w) \leq \alpha f(u) + (1-\alpha)f(w) - \frac{\lambda}{2}\alpha(1-\alpha)\|w-u\|^2 \quad (15)$$

$$\Rightarrow f(w + \alpha(u-w)) \leq \alpha(f(u) - f(w)) + f(w) - \frac{\lambda}{2}\alpha(1-\alpha)\|w-u\|^2 \quad (16)$$

$$\Rightarrow \frac{f(w + \alpha(u-w)) - f(w)}{\alpha} \leq f(u) - f(w) - \frac{\lambda}{2}(1-\alpha)\|w-u\|^2 \quad (17)$$

$$\Rightarrow \frac{f(w + \alpha(u-w)) - f(w)}{\alpha(u-w)}(u-w) \leq f(u) - f(w) - \frac{\lambda}{2}(1-\alpha)\|w-u\|^2 \quad (18)$$

Taking the limit as  $\alpha$  goes to 0, we have

$$\Rightarrow f'(w)(u-w) \leq f(u) - f(w) - \frac{\lambda}{2}\|w-u\|^2 \quad (19)$$

$$\Rightarrow \langle v, (u-w) \rangle \leq f(u) - f(w) - \frac{\lambda}{2}\|w-u\|^2 \quad (20)$$

Multiplying both sides by -1.

$$\Rightarrow \langle v, (w-u) \rangle \geq f(w) - f(u) + \frac{\lambda}{2}\|w-u\|^2 \quad (21)$$

Since  $\langle u, w \rangle = \langle w, u \rangle$  we have,

$$\langle (w-u), v \rangle \geq f(w) - f(u) + \frac{\lambda}{2}\|w-u\|^2 \quad (22)$$

### 1.3 Kernel construction

#### 1.3.1 a

Since

$$\alpha K_1(u, v) = \langle \sqrt{\alpha} \phi_1(u), \sqrt{\alpha} \phi_1(v) \rangle \quad (23)$$

and

$$\beta K_2(u, v) = \langle \sqrt{\beta} \phi_2(u), \sqrt{\beta} \phi_2(v) \rangle \quad (24)$$

Therefore

$$\begin{aligned} K(u, v) &= \alpha K_1(u, v) + \beta K_2(u, v) \\ &= \langle \sqrt{\alpha} \phi_1(u), \sqrt{\alpha} \phi_1(v) \rangle + \langle \sqrt{\beta} \phi_2(u), \sqrt{\beta} \phi_2(v) \rangle \\ &= \langle [\sqrt{\alpha} \phi_1(u), \sqrt{\beta} \phi_2(u)], [\sqrt{\alpha} \phi_1(v), \sqrt{\beta} \phi_2(v)] \rangle \end{aligned} \quad (25)$$

As we can see it is represented as an inner product of 2 vectors which can be interpreted as

$$\langle \Phi_1(u), \Phi_2(v) \rangle \quad (26)$$

#### 1.3.2 b

$$K_1(u, v) = \sum_i \phi_i(u), \phi_i(v) \quad (27)$$

$$K_2(u, v) = \sum_j \psi_j(u), \psi_j(v) \quad (28)$$

$$\begin{aligned} \implies K_1(u, v) K_2(u, v) &= \left( \sum_i \phi_i(u), \phi_i(v) \right) \left( \sum_j \psi_j(u), \psi_j(v) \right) \\ &= \sum_{i,j} \phi_i(u), \psi_j(u), \phi_i(v), \psi_j(v) \end{aligned} \quad (29)$$

We can write  $\Phi_k = \phi_i(u), \psi_j(u)$ .

$$\implies K_1(u, v) K_2(u, v) = \sum_k \Phi_k(u) \Phi_k(v) = K(u, v) \quad (30)$$

Therefore, product of valid kernels is a valid kernel.

### 1.4 Local minimum

Lets consider a sample point (x,y) such that  $x = c$  and  $y = 1$ . By just using a point in the sample space we can show that 0-1 loss function suffers from local minima.

Let  $w = -c$ ,  $\text{sign}(\langle w, x \rangle) = -1 \neq y$ , so  $L_s(w) = 1$ .  
Let  $\epsilon$  be a small scalar value and for every  $w'$ ,  $\|w' - w\| \leq \epsilon$ .

$$\begin{aligned}
\langle w', x \rangle &= \langle w, x \rangle + \langle w' - w, x \rangle \\
&= -c^2 + \langle w' - w, x \rangle \\
&\leq -c^2 + \|w' - w\| \|x\| \quad (\text{Using Cauchy Schwartz Inequality}) \\
&= -c^2 + \epsilon c \\
&< 0
\end{aligned} \tag{31}$$

Therefore  $L_s(w') = L_s(w) = 1$ , and hence  $w$  is a local minimum.  
Let  $w^* = c$ ,  $\langle w^*, x \rangle = c^2$  and so  $L_s(w^*) = 0$ .  
 $L_s(w^*) < L_s(w)$  which shows that  $w$  is not a global minimum.

## 1.5 Learnability of logistic regression

Let  $f(x) = \log(1 + e^x)$ . The gradient of  $f(x)$  is

$$f'(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}} < 1$$

for  $x \in R$ , so the function  $f(x)$  is 1-Lipschitz.

$g(-y\langle w, x \rangle)$  is B-Lipschitz since the norm is bounded by B.

Using claim 12.7 in the book, which says that if  $f(x) = g_1(g_2(x))$  and  $g_1$  is  $\rho_1$ -Lipschitz and  $g_2$  is  $\rho_2$ -Lipschitz, then  $f(x)$  is  $\rho_1\rho_2$ -Lipschitz.

Therefore  $l(w, \langle x, y \rangle)$  is  $(1*B)$ -Lipschitz = B-Lipschitz.

$$f''(x) = f' \left( \frac{1}{1 + e^{-x}} \right) = \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1}{(1 + e^x)(1 + e^{-x})} = \frac{1}{(2 + e^x + e^{-x})} \leq \frac{1}{4} \tag{32}$$

Since  $f''(x) \leq \frac{1}{4}$ , we can conclude that  $f'(x)$  is  $\frac{1}{4}$ -Lipschitz. Since  $f''(x)$  is nonnegative,  $f(x)$  is also convex. Therefore from the definition of smooth function we have that  $f(x)$  is a  $\frac{1}{4}$  smooth function. We know that composition of a smooth scalar function over a linear function preserves smoothness. Since  $l(w, \langle x, y \rangle) = f(-y\langle w, x \rangle)$ ,  $l$  is  $\frac{1}{4}\|x\|^2 = \frac{1}{4}B^2$  smooth.

The norm of the hypothesis class  $H$  is bounded above by B. Therefore it satisfies the first requirement of the Convex-Lipschitz Bounded and Convex-Smooth Bounded problem.

According to claim 12.4 in the book, the loss function  $l(w, \langle x, y \rangle)$  is convex since its a composition of the convex function  $f$  onto a linear function. The function is also B-Lipschitz as proved earlier. Therefore the problem is a Convex-Lipschitz Bounded problem with parameters  $B, B$ .

The loss function  $l(w, \langle x, y \rangle) = \log(1 + e^{-y\langle w, x \rangle})$  is non-negative since logarithm function is non-negative for values greater than 1. It is also convex and  $\frac{1}{4}B^2$  smooth as proved earlier. Therefore the problem is a Convex-Smooth Bounded problem with parameters  $\frac{1}{4}B^2, B$ .

## 1.6 Learnability of Halfspaces with hinge loss

$$l = \max\{0, 1 - y\langle w, x \rangle\}$$

For some  $x \in R^d$  and  $y \in \{-1, +1\}$ ,

$$g(w) = 1 - y\langle w, x \rangle \quad (33)$$

is a convex function since  $g'(w)$  is constant(monotonically non-decreasing) and  $g''(w) = 0$  (which is non-negative).

Using the property that the maximum of convex functions is also convex(from the Claim 12.5 in the book), we get that  $l = \max\{0, 1 - y\langle w, x \rangle\}$  is also a convex function.

Let  $w_1, w_2$  be two vectors such that  $w_1, w_2 \in R^d$ , then if we show that  $\|l_1 - l_2\| \leq R\|w_1 - w_2\|$  then we can say that  $l$  is  $R$ -lipschitz.

Case 1:  $(y\langle w_1, x \rangle \geq 1, y\langle w_2, x \rangle \geq 1)$

$$\begin{aligned} \|l_1 - l_2\| &= 0 - 0 \\ &\leq R\|w_1 - w_2\| \end{aligned} \quad (34)$$

Case 2:  $(y\langle w_1, x \rangle < 1, y\langle w_2, x \rangle \geq 1)$

$$\begin{aligned} \|l_1 - l_2\| &= l_1 - l_2 \\ &= (1 - y\langle w_1, x \rangle) - 0 \\ &< (1 - y\langle w_1, x \rangle) - (1 - y\langle w_2, x \rangle) \\ &= y\langle w_2 - w_1, x \rangle \\ &\leq \|w_2 - w_1\| \|x\| \text{ (Using Cauchy-Schwartz Inequality)} \\ &\leq R\|w_2 - w_1\| \end{aligned} \quad (35)$$

Case 3:  $(y\langle w_1, x \rangle \geq 1, y\langle w_2, x \rangle < 1)$

$$\begin{aligned} \|l_1 - l_2\| &= -(l_1 - l_2) \\ &= -(0 - (1 - y\langle w_2, x \rangle)) \\ &= (1 - y\langle w_2, x \rangle) \\ &< (1 - y\langle w_2, x \rangle) - (1 - y\langle w_1, x \rangle) \\ &= y\langle w_1 - w_2, x \rangle \\ &\leq \|w_2 - w_1\| \|x\| \text{ (Using Cauchy-Schwartz Inequality)} \\ &\leq R\|w_2 - w_1\| \end{aligned} \quad (36)$$

Case 4:  $(y\langle w_1, x \rangle < 1, y\langle w_2, x \rangle < 1)$

Lets suppose  $(1 - y\langle w_1, x \rangle) \geq (1 - y\langle w_2, x \rangle)$ .

$$\begin{aligned} \|l_1 - l_2\| &= (l_1 - l_2) \\ &= (1 - y\langle w_1, x \rangle) - (1 - y\langle w_2, x \rangle) \\ &= y\langle w_2 - w_1, x \rangle \\ &\leq \|w_2 - w_1\| \|x\| \text{ (Using Cauchy-Schwartz Inequality)} \\ &\leq R\|w_2 - w_1\| \end{aligned} \quad (37)$$

The same is applicable when  $(1 - y\langle w_1, x \rangle) \leq (1 - y\langle w_2, x \rangle)$ . Therefore the function  $l$  is  $R$ -lipschitz.