

Assignment 2

Manasa Navalgund(112871754)

April 20, 2020

1 Theory

Q1.

$$\begin{aligned} \text{Given } N(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\} \\ E[x] &= \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) x dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\} x dx \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\} x dx \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp\left(-\left(\frac{x - \mu}{\sqrt{2}\sigma}\right)^2\right) x dx \end{aligned}$$

$$\begin{aligned} \text{Substituting } t &= \frac{x - \mu}{\sqrt{2}\sigma} \\ &= dt = \frac{dx}{\sqrt{2}\sigma} \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp(-t^2) (\sqrt{2}\sigma t + \mu)(\sqrt{2}\sigma dt) \\ &= \frac{\sqrt{2}\sigma}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp(-t^2) (\sqrt{2}\sigma t + \mu) dt \\ &= \frac{1}{\sqrt{\pi}} \left(\int_{-\infty}^{\infty} \sqrt{2}\sigma t \exp(-t^2) dt + \int_{-\infty}^{\infty} \exp(-t^2) \mu dt \right) \end{aligned}$$

From Gaussian integral, [1]

$$\int_{-\infty}^{\infty} \exp(-t^2) dt = \sqrt{\pi}$$

Substituting this in above equation,

$$\begin{aligned} &= \frac{1}{\sqrt{\pi}} \left(\int_{-\infty}^{\infty} \sqrt{2}\sigma t \exp(-t^2) dt + \mu\sqrt{\pi} \right) \\ \frac{d}{dt} \exp(-t^2) &= -2t \exp(-t^2) \end{aligned}$$

Susbtituting this in above,

$$\begin{aligned} &= \frac{1}{\sqrt{\pi}} \left(\sqrt{2}\sigma \left[\frac{-1}{2} \exp(-t^2) \right]_{-\infty}^{\infty} + \mu\sqrt{\pi} \right) \\ &= \frac{\mu\sqrt{\pi}}{\sqrt{\pi}} \\ &= \mu \end{aligned}$$

Therefore $E[x] = \mu$

Hence proved.

We know that,

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\}$$

And since this is a valid probability distribution,

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\} dx = 1$$

Differentiating this w.r.t σ^2

$$\frac{d}{d\sigma^2} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\} dx = \frac{d}{d\sigma^2}(1)$$

Using Lebiniz theorem, we can interchange integration and differentiation[\[3\]](#),

$$\int_{-\infty}^{\infty} \frac{d}{d\sigma^2} \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\} \right) dx = \frac{d}{d\sigma^2}(1)$$

Substituting $\sigma^2 = t$

$$\int_{-\infty}^{\infty} \frac{d}{dt} \left(\frac{1}{\sqrt{2\pi t}} \exp\left\{\frac{-1}{2t}(x - \mu)^2\right\} \right) dx = \frac{d}{dt}(1)$$

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

Substituting $t = \sigma^2$

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \left(\frac{1}{\sqrt{\sigma^2}} \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\} \left(\frac{-(x - \mu)^2}{2} \left(\frac{-1}{\sigma^4} \right) \right) + \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\} \left(\frac{-1}{2\sigma^3} \right) \right) dx = 0$$

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{\frac{-1}{2\sigma^2}(x - \mu)^2\right\} \left(\left(\frac{(x - \mu)^2}{2\sigma^4} \right) + \left(\frac{-1}{2\sigma^2} \right) \right) dx = 0$$

$$\int_{-\infty}^{\infty} N \left(\left(\frac{(x - \mu)^2}{2\sigma^4} \right) - \left(\frac{1}{2\sigma^2} \right) \right) dx = 0$$

$$\begin{aligned} \int_{-\infty}^{\infty} N \left((x - \mu)^2 - \sigma^2 \right) dx &= 0 \\ \int_{-\infty}^{\infty} N \left(x^2 + \mu^2 - 2x\mu - \sigma^2 \right) dx &= 0 \end{aligned}$$

$$\int_{-\infty}^{\infty} N(x^2) dx + \int_{-\infty}^{\infty} N(\mu^2) dx - \int_{-\infty}^{\infty} N(2x\mu) dx - \int_{-\infty}^{\infty} N(\sigma^2) dx = 0$$

Since μ and σ are constants w.r.t x and $\int_{-\infty}^{\infty} N dx = 1$ and above result,

$$E[x^2] + \mu^2 - 2\mu^2 - \sigma^2 = 0$$

$$E[x^2] = \mu^2 + \sigma^2$$

Hence proved.

$$\text{var}(x) = E[x^2] - E[x]^2$$

From above results,

$$\text{var}(x) = \mu^2 + \sigma^2 - (\mu)^2$$

$$\text{var}(x) = \sigma^2$$

Q2.

Given: f is a strongly-convex function with parameter λ

For a strongly-convex function, we have, DEFINITION 13.4, UML

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

Dividing the definition by α , we get,

$$\frac{f(w + \alpha(u - w)) - f(w)}{\alpha} \leq f(u) - f(w) - \frac{\lambda}{2}(1 - \alpha)\|w - u\|^2 \quad (3)$$

We know from LEMMA 14.3 and DEFINITION 14.4 of UML that a subgradient v , satisfies following equation,

$$f(u) - f(w) \geq \langle u - w, v \rangle$$

Substituting $u = w + \alpha(u - w)$,

$$f(w + \alpha(u - w)) - f(w) \geq \langle w + \alpha(u - w) - w, v \rangle$$

$$f(w + \alpha(u - w)) - f(w) \geq \alpha \langle u - w, v \rangle$$

Substituting this in (3) above,

$$\langle u - w, v \rangle \leq f(u) - f(w) - \frac{\lambda}{2}(1 - \alpha)\|w - u\|^2$$

Multiplying both sides by -1, we get:

$$\implies \langle w - u, v \rangle \geq f(w) - f(u) + \frac{\lambda}{2}(1 - \alpha)\|w - u\|^2$$

For $\alpha > 0$, Since u is minimizer of f , it follows $\alpha = 0$ is minimizer

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

Q3.

Given: K_1 and K_2 are valid kernels over the domain of X .

a. $K(u, v) = \alpha K_1(u, v) + \beta K_2(u, v)$

$\alpha K_1(u, v)$ can be represented as : $\langle \sqrt{\alpha} \Psi_1(u), \sqrt{\alpha} \Psi_1(v) \rangle$

$\beta K_2(u, v)$ can be represented as : $\langle \sqrt{\beta} \Psi_2(u), \sqrt{\beta} \Psi_2(v) \rangle$

$$\begin{aligned} \text{Now, } K(u, v) &= \alpha K_1(u, v) + \beta K_2(u, v) \\ &= \langle \sqrt{\alpha} \Psi_1(u), \sqrt{\alpha} \Psi_1(v) \rangle + \langle \sqrt{\beta} \Psi_2(u), \sqrt{\beta} \Psi_2(v) \rangle \\ &= \langle \sqrt{\alpha} \Psi_1(u) \sqrt{\beta} \Psi_2(u), \sqrt{\alpha} \Psi_1(v) \sqrt{\beta} \Psi_2(v) \rangle \end{aligned}$$

Since, α, β are positive, $\sqrt{\alpha}$ and $\sqrt{\beta}$ are positive scalars and we know Ψ_1, Ψ_2 are valid kernels,

Therefore this shows above inner product is a valid kernel

b. $K(u, v) = K_1(u, v) K_2(u, v)$

The gram matrix of this product is given as $K = K_1 \cdot K_2$

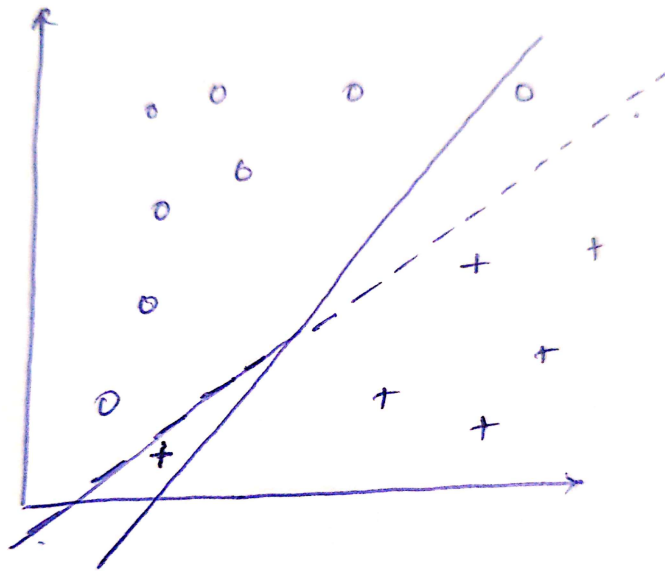
where \cdot represents Hadamard product.

Since K_1 and K_2 are positive and symmetric positive definite matrices. When we perform hadamard product on eigen decompositions, it gives a symmetric positive definite matrix as result. Since, this matrix suffices conditions for gram matrix of some two arguments, it is valid kernel function.

Q4.

Consider an example Halfspaces with loss function as $l = \max\{0, 1 - \langle w, x \rangle\}$

Consider the following training data set,



Scanned with CamScanner

Figure 1: Training data

Here, with very small step size, the algorithm reaches till solid line and number of epochs end before line converges to dashed line. Since this has minimum loss till that point, algorithm returns dashed line as best hyperplane separating the data.

However, with more epochs or larger step size, hyperplane could have converged to dashed line which classifies the data perfectly. This means solid line is global minima but algorithm converged at solid line which is local minima.

Therefore, halfspace class of functions can suffer from local minima

Q5

Convex-Lipschitz bounds

Given function $l = g = \log(1 + \exp(-y < w, x >))$

$$g'(x) = \left| \frac{\exp(-y < w, x >)}{1 + \exp(-y < w, x >)} \right|$$

$$= \left| \frac{1}{1 + \exp(y < w, x >)} \right|$$

Here, we can see that $g'(x)$ is monotonically increasing function i.e., g'' is non negative. Therefore, l is a convex function and given that $\|x\| \leq B$

We know from DEFINITION 12.6 UML, that Lipschitzness is defined as:

$$\|f(w_1) - f(w_2)\| \leq \rho \|w_1 - w_2\|$$

By mean value theorem,

$$\|f(w_1) - f(w_2)\| \leq f'(u) \|w_1 - w_2\|$$

Given loss function is $g = \log(1 + \exp(-y < w, x >))$

$$g'(x) = \left| \frac{\exp(-y < w, x >)}{1 + \exp(-y < w, x >)} \right|$$

$$= \left| \frac{1}{1 + \exp(y < w, x >)} \right| \leq 1$$

Therefore, l is 1-Lipschitz.

But from claim 12.7, Lipschitzness for l is $\rho_1 \|w\| = B$

l is B -Lipschitz

Therefore, l is convex-Lipschitz bounded with parameters $B, \rho=B$

Convex-Smooth-bounded

Given function $l = g = \log(1 + \exp(-y < w, x >))$

From the definition of convex-smooth-bounded, we know that a learning problem is convex-smooth-bounded if,

1. For all $w \in H, \|w\| \leq B$

2. loss function is convex, non-negative and β -smooth

$$g'(x) = \left| \frac{\exp(-y < w, x >)}{1 + \exp(-y < w, x >)} \right|$$

$$= \left| \frac{1}{1 + \exp(y < w, x >)} \right|$$

Here, we can see that $g'(x)$ is monotonically increasing function i.e., g'' is non negative. Therefore, l is a convex function and given that $\|x\| \leq B$

loss function = $g = \log(1 + \exp(-y < w, x >))$

For a function to be β -smooth, its gradient has to be β -Lipschitz (From DEFINITION 12.8 UML),

$$g'(x) = \left| \frac{\exp(-y < w, x >)}{1 + \exp(-y < w, x >)} \right|$$

$$g''(x) = \left| \frac{\exp(-y < w, x >)}{(1 + \exp(-y < w, x >))^2} \right|$$

$$g''(x) = \left| \frac{1}{(1 + \exp(-y < w, x >))(1 + \exp(y < w, x >))} \right| \leq \frac{1}{4}$$

Therefore, gradient is $\frac{1}{4}$ -Lipschitz and loss function is $\frac{1}{4}$ -smooth.

But from claim 12.9, smoothness for l is $\beta ||x||^2 = \frac{B^2}{4}$.

l is $\frac{B^2}{4}$ -smooth There we have proven convex-smooth-bounded with parameters $B, \frac{B^2}{4}$.

Q6. Given $l(w, (x, y)) = \max\{0, 1 - y < w, x >\}$

Proof of convexity:

Proving l to be R-Lipschitz:

We know that Lipschitzness is given by :

$$||l(w_1) - l(w_2)|| \leq \rho ||w_1 - w_2||$$

Let's consider a point (x,y) such that:

$$l_{w_1} = \max\{0, 1 - y < w_1, x >\}$$

$$l_{w_2} = \max\{0, 1 - y < w_2, x >\}$$

There can be multiple cases:

Case 1: $l_{w_1} = l_{w_2} = 0$ In this case, from definition of Lipschitzness, we have,

$$|l_{w_1} - l_{w_2}| = 0 \leq R|w_1 - w_2| \text{ Since } R \geq 0 \text{ and } |w_1 - w_2| \geq 0$$

Case 2: When atleast one of l_{w_1} or l_{w_2} is not 0

$$|l_{w_1} - l_{w_2}| = |(1 - \max\{0, 1 - y < w_1, x >\}) - (1 - \max\{0, 1 - y < w_2, x >\})|$$

Lets assume that $l_{w_1} \geq l_{w_2}$

$$\begin{aligned} |l_{w_1} - l_{w_2}| &= l_{w_1} - l_{w_2} \\ &= (1 - y < w_1, x >) - (\max\{0, 1 - y < w_2, x >\}) \\ &\leq (1 - y < w_1, x >) - (1 - y < w_2, x >) = y < w_2 - w_1, x > \\ &\leq ||w_2 - w_1|| |x| \\ &\leq ||w_1 - w_2|| \end{aligned}$$

Similarly if we assume $l_{w_2} \geq l_{w_1}$,

$$\begin{aligned} |l_{w_1} - l_{w_2}| &= l_{w_2} - l_{w_1} \\ &= (1 - y < w_2, x >) - (\max\{0, 1 - y < w_1, x >\}) \\ &\leq (1 - y < w_2, x >) - (1 - y < w_1, x >) = y < w_1 - w_2, x > \\ &\leq ||w_1 - w_2|| |x| \\ &\leq ||w_1 - w_2|| \end{aligned}$$

Therefore, $l(w, (x, y)) = \max\{0, 1 - y < w, x >\}$ is R-Lipschitz

2 Programming

I get the following plot for the data:

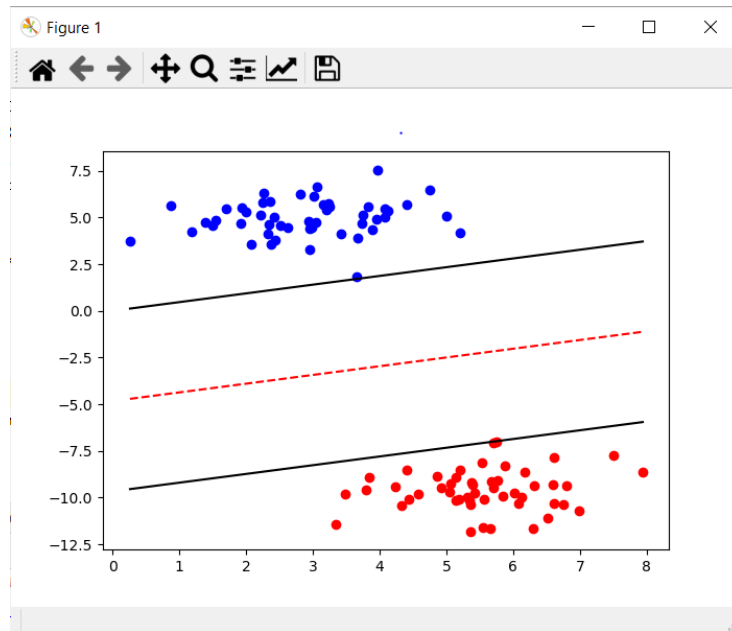


Figure 2: Training data

Output is redirected to file 'ReadMe.txt'

Output of above run:

the total number of data points on which your test method was run, : 20

the total number of data points on which your test method was misclassified, :

0

References

1. https://en.wikipedia.org/wiki/Gaussian_integral
2. Understanding Machine Learning: From Theory to Algorithms
3. <https://math.stackexchange.com/questions/2530213/when-can-we-interchange-integration-and-differentiation>