problem 1

a) $w_i \sim N(0, T^2)$

$$p(w \mid 0, T^2) = \prod_{i=1}^{d} \frac{1}{T\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{w_i}{T}\right)^2}$$

$$w_{MAP} = \arg\max_{w} \left( \sum_{i=1}^{N} \ln p(y^{(i)} \mid x^{(i)}, w) + \ln p(w) \right)$$

$$f = \sum_{i=1}^{N} \ln \frac{1}{\sqrt{2\pi}\sigma} \exp\left( \frac{-(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2} \right) + \sum_{i=1}^{d} \ln\left( \frac{1}{T\sqrt{2\pi}} \exp\left( -\frac{1}{2}\left(\frac{w_i}{T}\right)^2 \right) \right)$$

$$= \ln \frac{n}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^{N} (y^{(i)} - w^T x^{(i)})^2 + \ln \frac{d}{T\sqrt{2\pi}} - \frac{1}{2T^2} \sum_{i=1}^{d}(w_i^2)$$

$$= \ln \frac{n}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2}(y - xw)^T (y - xw) + \ln \frac{d}{T\sqrt{2\pi}} - \frac{1}{2T^2} w^T w$$

Take derivative

$$\frac{df}{dw} = \frac{1}{\sigma^2}(x^T y - x^T x w) - \frac{1}{T^2} w$$

$$= \frac{1}{\sigma^2} x^T(y - xw) - \frac{1}{T^2} w \overset{set}{=} 0$$

$$\frac{1}{\sigma^2} x^T(y - xw) = \frac{1}{T^2} w$$

$$x^T y - x^T x w = \frac{\sigma^2}{T^2} w$$

$$w_{MAP} = x^T y \cdot \left( x^T x + \frac{\sigma^2}{T^2} I \right)^{-1}$$

From ridge regression, we have similar form

$$w_{RR} = (\lambda I + x^T x)^{-1} \cdot x^T y$$

b) $w^{(i)} \sim Laplace(0, b)$

$$p(w_i \mid b) = \prod_{i=1}^{d} \frac{1}{2b} e^{-\frac{|w_i|}{b}}$$

$$w_{MAP} = \arg\max \left( \sum_{i=1}^{N} \ln \frac{1}{\sqrt{2\pi}\sigma} \exp\left( \frac{-(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2} \right) + \sum_{i=1}^{d} \ln \frac{1}{2b} \exp\left( -\frac{|w_i|}{b} \right) \right)$$

$$= \arg\max \left( -\frac{1}{2\sigma^2} \sum_{i=1}^{N} (y^{(i)} - w^T x^{(i)})^2 - \frac{1}{b} \sum_{i=1}^{d} |w_i| \right)$$

$$= \arg\min \left( \sum_{i=1}^{N} (y^{(i)} - w^T x^{(i)})^2 + \frac{2\sigma^2}{b} \|w\| \right)$$

From Lasso regression, we see similar formula

$$\hat{w}_{Lasso} = \arg\min \frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \hat{w}^T x^{(i)})^2 + \lambda \|w\|$$

Problem 2.

(1)  From bayes rule :

$$f_{bayes}(x) = \arg\max\ p(x|y) \cdot pr(y) \qquad \text{where}\quad p(x|y) \text{ is class conditional density}$$
$$\text{and}\quad pr(y) \text{ is class prior}$$

Assume we have data
$$\{(x^{(i)}, y^{(i)})\}_{i=1}^N = \{(x^{(1)}, y^{(1)}), \dots (x^{(N)}, y^{(N)})\}$$
$$x^{(i)} \in X \subseteq R^d \quad y^{(i)} \in Y (:= \{0,1\})$$

Then from Gaussian Assumption
$$P(x,y) = p(y)\,p(x|y) = \begin{cases} p_0 \frac{1}{\sqrt{2\pi}\,\delta_0}\,e^{\frac{(x-u_0)^2}{2\delta_0^2}} & \text{if } y=0 \\ p_1 \frac{1}{\sqrt{2\pi}\,\delta_1}\,e^{-\frac{(x-u_0)^2}{2\delta_1}} & \text{if } y=1 \end{cases}$$

The Bayes optimal one under the assumed joint distribution depends on

$$\mathbb{1}\left(Pr(y=1|x) \ge Pr(y=0|x)\right) \overset{bayes}{\Rightarrow} \mathbb{1}\left(P(x|y=1)\,Pr(y=1) \ge P(x|y=0)\,Pr(y=0)\right)$$

$$\Rightarrow \mathbb{1}\left(-\frac{(x-u_1)^2}{2\delta_1^2} - \log\sqrt{2\pi}\,\delta_1 + \log p_1 \ge -\frac{(x-u_0)^2}{2\delta_0^2} - \log\sqrt{2\pi}\,\delta_0 + \log p_0\right)$$

$$\Rightarrow \mathbb{1}(ax^2 + bx + c \ge 0) \qquad \text{So the decision bondary is not linear}$$

In matrix form

Gaussian Assumption
$$P(x|y, u_y, \Sigma_y) = \frac{1}{(2\pi)^{d/2}\,|\Sigma_y|^{1/2}}\,\exp\left(\frac{(x-u_y)^T(\Sigma_y)^{-1}(x-u_y)}{2}\right)$$

Let $d_\Sigma(x,u) = (x-u)^T(\Sigma)^{-1}(x-u)$

$$f(x) = \mathbb{1}\left(Pr(y=1|x) > Pr(y=0|x)\right) \Rightarrow \mathbb{1}\left(\ln\frac{Pr(y=1|x)}{Pr(y=0|x)} > 0\right)$$

$$\Rightarrow \mathbb{1}\left(\ln\frac{P(x|y,u,\Sigma)\,Pr(y=1)}{P(x|y,u,\Sigma)\,Pr(y=0)}\right) \qquad * (Pr(y=1) = \pi_1)$$

$$\Rightarrow \mathbb{1}\left(\ln\frac{\pi_1}{1-\pi_1} - \frac{1}{2}\ln\frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2}(d_{\Sigma_1}(x,u_1) - d_{\Sigma_0}(x,u_0)) > 0\right)$$

Since $d_\Sigma(x,u)$ is in quadratic form, GDA has a quadratic decision boundary.

Further, if we assume $\Sigma_1 = \Sigma_0 = \Sigma$
$$d_{\Sigma_1}(x,u_1) - d_{\Sigma_0}(x,u_0) = (x-u_1)^T(\Sigma)^{-1}(x-u_1) - (x-u_0)^T(\Sigma)^{-1}(x-u_0)$$
$$= x^T\Sigma^{-1}x - x^T\Sigma^{-1}u_1 - u_1^T\Sigma^{-1}x + u_1^Tu_1\Sigma^{-1}$$
$$- x^T\Sigma^{-1}x + x^T\Sigma^{-1}u_0 + u_0^T\Sigma^{-1}x - u_0^T\Sigma^{-1}u_0$$

$$= x^T\Sigma^{-1}(u_0-u_1) - u_1^T\Sigma^{-1}u_1 + u_2^T\Sigma^{-1}u_2$$

So it is linear in X.

(2)  From   Gaussian Discriminant Analysis

$$f(x) = \underset{y \in \{0,1\}}{\arg\max} \; Pr(y|x) = \underset{y \in \{0,1\}}{\arg\max} \; \underbrace{Pr(y)}_{\pi_y} P(x|y; \mu_y, \Sigma_y)$$

$$= \underset{y}{\arg\max} \; \frac{1}{(2\pi)^{d/2} |\Sigma_y|^{1/2}} \exp\left\{ -\frac{1}{2}(x-\mu_y)^T \Sigma_y^{-1}(x-\mu_y)\right\} \cdot \pi_y$$

$$= \mathbb{1}\left[ Pr(y=1|x) > Pr(y=0|x)\right]$$

$$= \mathbb{1}\left[ \ln\frac{\pi_1}{1-\pi_1} - \frac{1}{2}\ln\frac{\Sigma_1}{\Sigma_0} - \frac{1}{2}\left(d_{\Sigma_1}(x,\mu_1) - d_{\Sigma_0}(x,\mu_0)\right) > 0\right]$$

Since $\Sigma_1, \Sigma_0$ are diagonal.   $|\Sigma_1| = \prod_{k=1}^{d} \sigma_k$, $y=1$   $|\Sigma_0| = \prod_{k=1}^{d} \sigma_k$, $y=0$

$$d_{\Sigma_y}(x, \mu) = \underset{1 \times n}{(x-\mu)^T} \; \underset{n \times n}{(\Sigma)^{-1}} \; \underset{n \times 1}{(x-\mu)} = \sum_{k=1}^{d}\left(\frac{x-\mu}{\sigma}\right)^2$$

$$= \mathbb{1}\left[ \ln\frac{\pi_1}{1-\pi_1} - \frac{1}{2}\sum_{k=1}^{d}\ln\frac{\sigma_{k,y=1}}{\sigma_{k,y=0}} - \frac{1}{2}\left( \sum_{k=1}^{d}\left(\frac{x_k - \mu_{k,y=1}}{\sigma_{k,y=1}}\right)^2 - \sum_{k=1}^{d}\left(\frac{x_k - \mu_{k,y=0}}{\sigma_{k,y=0}}\right)^2\right) > 0\right]$$

From   Gaussian Naive Bayes classifier

$$f(x) = \underset{y \in y}{\arg\max} \prod_{k=1}^{d} P(x_k|y; \mu_k, \sigma_k^2) \cdot \underbrace{Pr(y)}_{\pi_y} \qquad x_k \in \{0,1\}$$

$$= \mathbb{1}\left[ \prod_{k=1}^{d}\frac{1}{\sigma\sqrt{2\pi}}\exp\left\{-\frac{1}{2}\left(\frac{x_k - \mu_{k,y=1}}{\sigma_{k,y=1}}\right)^2\right\} \cdot \pi_1 > \prod_{k=1}^{d}\frac{1}{\sigma\sqrt{2\pi}}\exp\left\{-\frac{1}{2}\left(\frac{x_k - \mu_{k,y=0}}{\sigma_{k,y=0}}\right)^2\right\} \cdot \pi_0\right]$$

$$= \mathbb{1}\left[ \ln \frac{\prod_{k=1}^{d}\frac{1}{\sigma_k\sqrt{2\pi}}\exp\left\{\frac{1}{2}\left(\frac{x_k - \mu_{y=1}}{\sigma_{y=1}}\right)^2\right\}\pi_1}{\prod_{k=1}^{d}\frac{1}{\sigma_k\sqrt{2\pi}}\exp\left\{\frac{1}{2}\left(\frac{x_k - \mu_{y=0}}{\sigma_{k,y=0}}\right)^2\right\}\pi_2} > 0 \right]$$

$$= \mathbb{1}\left[ \ln\frac{\pi_1}{1-\pi_1} - \frac{1}{2}\sum_{k=1}^{d}\ln\frac{\sigma_{k,y=0}}{\sigma_{k,y=1}} - \frac{1}{2}\left( \sum_{k=1}^{d}\left(\frac{x_k - \mu_{k,y=0}}{\sigma_{k,y=0}}\right)^2 - \sum_{k=1}^{d}\left(\frac{x_k - \mu_{k,y=0}}{\sigma_{k,y=0}}\right)^2\right) > 0\right]$$

Therefore, when the covariance matrix $\Sigma_y$ is diagonal
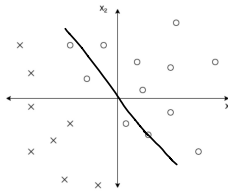Gaussian NBC is a special case of Gaussian Discriminant Analysis.

problem 3

a) optimization function :

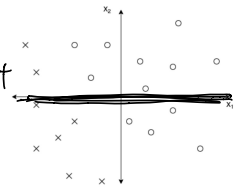$$\sum_{i=1}^{N} \log(P(y^{(i)} \mid x^{(i)}; w_0, w_1, w_2)) - \lambda w_j^2$$

when $\lambda$ is very large, the penalty term has large impact. To minimize the cost function, $w_j$ has to be 0.

when $w_0 = 0$,
train error increase,
arround 3-4 points
misclassified.

when $w_1 = 0$
train error increase most
arround 6 points
misclassified.

when $w_2 = 0$
train error increase
arround 3 points
misclassified.