## HOMEWORK-3

**Q1.** Given,

$$\min_{\theta, b, \xi_i} \|\underline{\theta}\|_2^2 + C \sum_{i=1}^{n} \xi_i^2 \qquad s.t. \quad y_i(\underline{\theta}^T \underline{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0 \quad \forall i$$

We define the following Lagrangian:

$$\mathcal{L}(\underline{\theta}, b, \xi, \underline{\alpha}, \underline{r}) = \|\underline{\theta}\|^2 + C \sum_{i=1}^{N} \xi_i^2 + \sum_{i=1}^{N} \alpha_i (1 - \xi_i - y_i(\underline{\theta}^T \underline{x}_i + b)) - \sum_{i=1}^{N} r_i \xi_i$$

$$\frac{\partial \mathcal{L}}{\partial \underline{\theta}} = 0 \Rightarrow 2\underline{\theta} - \sum_{i=1}^{N} \alpha_i y_i \underline{x}_i = 0 \Rightarrow \underline{\theta} = \frac{1}{2} \sum_{i=1}^{N} \alpha_i y_i \underline{x}_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow -\sum_{i=1}^{N} \alpha_i y_i = 0 \Rightarrow \sum_{i=1}^{N} \alpha_i y_i = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = 0 \Rightarrow 2C\xi_i - \alpha_i - r_i = 0 \Rightarrow \xi_i = \frac{1}{2C}\left(\alpha_i + r_i\right)$$

Now,

$$\mathcal{L} = \underline{\theta}^T \underline{\theta} + \frac{1}{4C} \sum_{i=1}^{N} (\alpha_i + r_i)^2 + \sum_{i=1}^{N} \alpha_i - \sum_{i=1}^{N} \alpha_i \xi_i - \underline{\theta}^T \sum_{i=1}^{N} \alpha_i y_i \underline{x}_i - b\sum_{i=1}^{N}\alpha_i y_i - \sum r_i \xi_i$$

$$= -\frac{1}{4} \sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \underline{x}_i^T \underline{x}_j + \frac{1}{4C} \sum_{i=1}^{N} (\alpha_i + r_i)^2 + \sum_{i=1}^{N} \alpha_i - \sum_{i=1}^{N} \frac{1}{2C}(\alpha_i + r_i)^2$$

$$\therefore G(\underline{\alpha}, \underline{r}) = -\frac{1}{4}\sum_{i=1}^{N}\sum_{j=1}^{N} y_i y_j \alpha_i \alpha_j \underline{x}_i^T \underline{x}_j + \sum_{i=1}^{N}\alpha_i - \frac{1}{4C}\sum_{i=1}^{N}(\alpha_i + r_i)^2$$

Thus, the dual SVM is the maximization (i.e. −minimization) of $G(\underline{\alpha}, \underline{r})$.

$$\therefore \min_{\underline{\alpha}, \underline{r}} \frac{1}{4}\sum\sum \alpha_i \alpha_j y_i y_j \underline{x}_i^T \underline{x}_j - \sum_{i=1}^{N}\alpha_i + \frac{1}{4C}\sum_{i=1}^{N}(\alpha_i + r_i)^2$$

$$s.t. \quad \sum_{i=1}^{N}\alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C \quad \forall \ i = 1, \dots, N$$

Q.3 We define model parameters as:

$$E[Z] = \mu \quad (\text{Bias})$$
$$E[(Z-\mu)^2] = \sigma^2 \quad (\text{Variance})$$
$$E\; Cor(Z_i, Z_j) = f_{ij} \quad (\text{Correlation})$$

Now,

For individual models,

$$\text{Bias} = E[Z_b] = \mu \quad \{b = 1, 2, \ldots, B\} \quad \{\because \text{Uniform sampling}\}$$
$$\text{Variance} = E[(Z_b-\mu)^2] = \sigma^2 \quad \{b = 1, 2, \ldots, B\}$$
$$\text{Correlation} = Cor(Z_i, Z_j)_{i \neq j} = f \quad \{i,j = 1, 2, \ldots B\}$$

Considering bagged models:

$$\text{Bias} = E\left[\frac{1}{B}\sum_{b=1}^{B} Z_b\right] = \frac{1}{B}\sum_{b=1}^{B} E[Z_b] = \frac{B \times \mu}{B} = \mu$$

⇒ Bias remains the same in both models.

$$\text{Variance} = E\left[\left(\frac{1}{B}\sum_{b=1}^{B} Z_b\right)^2\right] - \left(E\left[\sum_{b=1}^{B}\frac{1}{B}Z_b\right]\right)^2$$

$$E\left[\left(\frac{1}{B}\sum_{b=1}^{B} Z_b\right)^2\right] = \frac{1}{B^2} E\left[\sum_b Z_b^2 + \sum_{i \neq j, \, i,j=1,\ldots B} Z_i Z_j\right]$$

$$= \frac{1}{B^2}\sum_b E[Z_b^2] + \frac{1}{B^2} E\left[\sum Z_i Z_j\right]$$

$$= \frac{B(\sigma^2+\mu^2)}{B^2} + \frac{(B-1)B(Cov(Z_i, Z_j) + \mu^2)}{B^2}$$

$$= \frac{\sigma^2 + \mu^2}{B} + \frac{(B-1)(f\sigma^2 + \mu^2)}{B}$$

$$= \frac{\sigma^2(f(B-1)+1) + B\mu^2}{B}$$

$$\therefore \text{Variance} = \frac{\sigma^2(f(B-1)+1)}{B} + \mu^2 - \mu^2 = \frac{\sigma^2}{B} + f\sigma^2\left(1 - \frac{1}{B}\right)$$

Since $\rho \leq 1$

$\Rightarrow \rho\sigma^2\left(1-\dfrac{1}{B}\right) \leq \sigma^2\left(1-\dfrac{1}{B}\right)$

$\Rightarrow \rho\sigma^2\left(1-\dfrac{1}{B}\right) + \dfrac{\sigma^2}{B} \leq \sigma^2$

Clearly, the variance for bagging model is lesser than the individual model.

The bias - variance decomposition suggests:

$$E\left[(y-\hat{y})^2\right] = \text{Bias}^2 + \text{Variance} + \text{Irreducible error}$$

Thus, the bias is unchanged while variance reduces in bagging resulting in a net reduction of the squared error loss.

Q.4. Given,

$$\hat{\theta} = \arg\max_{\theta} \sum_{i=1}^{N}\sum_{m=1}^{M} w_m^{(i)}\left(\ln \mathcal{N}\left(x^{(i)} \mid \mu_m, \Sigma_m\right) + \ln \pi_m\right)$$

Let $J(\theta) = \sum_{i=1}^{N}\sum_{m=1}^{M} w_m^{(i)}\left(\ln\left\{\mathcal{N}\left(x^{(i)} \mid \mu_m, \Sigma_m\right)\right\} + \ln \pi_m\right)$ $\{\theta = \{\mu_m, \Sigma_m, \pi_m\}\}$

For maxima, $\dfrac{\partial J}{\partial \theta} = 0$

We know the estimated $\hat{\mu}_m$ for the current M-step, thus we do only need to update $\Sigma_m$.

$\therefore \dfrac{\partial J}{\partial \Sigma_m} = 0$

$J = \sum_{i=1}^{N}\sum_{m=1}^{M} w_m^{(i)}\left(\ln\left(\dfrac{1}{(2\pi)^{P/2}|\Sigma_m|^{1/2}} e^{-\frac{1}{2}(x-\mu_m)^T \Sigma_m^{-1}(x-\mu_m)}\right) + \ln \pi_m\right)$

$= \sum_{i=1}^{N}\sum_{m=1}^{M} w_m^{(i)}\left[-\dfrac{p}{2}\ln(2\pi) - \dfrac{1}{2}\ln(|\Sigma_m|) - \dfrac{1}{2}(x-\mu_m)^T \Sigma_m^{-1}(x-\mu_m) + \ln \pi_m\right]$

Let $\underline{\underline{\Lambda}}_m = \underline{\underline{\Sigma}}_m^{-1} \implies |\underline{\underline{\Lambda}}_m| = \dfrac{1}{|\underline{\underline{\Sigma}}_m|}$

$\therefore J = \sum\limits_{i=1}^{N} \sum\limits_{m=1}^{M} w_m^{(i)} \left( \dfrac{-p}{2} \ln(2\pi) + \dfrac{1}{2} \ln(|\underline{\underline{\Lambda}}_m|) - \dfrac{1}{2} (\underline{x} - \underline{\mu}_m)^T \underline{\underline{\Lambda}}_m (\underline{x} - \underline{\mu}_m) + \ln \pi_m \right)$

$\dfrac{\partial J}{\partial \underline{\underline{\Sigma}}_m} = \dfrac{\partial J}{\partial \underline{\underline{\Lambda}}_m} \dfrac{\partial \underline{\underline{\Lambda}}_m}{\partial \underline{\underline{\Sigma}}_m}$

$\implies \text{Maximi} \quad \dfrac{\partial J}{\partial \underline{\underline{\Lambda}}_m} = 0 \implies \dfrac{\partial J}{\partial \underline{\underline{\Sigma}}_m} = 0 \quad \text{since} \quad \dfrac{\partial \underline{\underline{\Lambda}}_m}{\partial \underline{\underline{\Sigma}}_m} \text{ is non-zero} \quad \{ \underline{\underline{\Sigma}}_m^{-1} \text{ exist}$

$\therefore \dfrac{\partial J}{\partial \underline{\underline{\Lambda}}_m} = \sum\limits_{i=1}^{N} \sum\limits_{m=1}^{M} w_m^{(i)} \left( \dfrac{1}{2} \dfrac{1}{|\underline{\underline{\Lambda}}_m|} \text{adj}(\underline{\underline{\Lambda}}_m) - \dfrac{1}{2} (\underline{x} - \underline{\mu}_m)(\underline{x} - \underline{\mu}_m)^T \right) = 0$

$\{ \text{Using identities:} \quad \dfrac{d|A|}{dA} = \dfrac{\text{adj}(A)}{|A|} \quad \& \quad \dfrac{d(\underline{x}^T A \underline{x})}{dA} = \underline{x}\underline{x}^T \}$

Also, we know that $\underline{\underline{\Sigma}}_m$ and thus $\underline{\underline{\Lambda}}_m$ are all equal.

$\implies 0 = \dfrac{1}{2} \dfrac{\text{adj}(\underline{\underline{\Lambda}}_m)}{|\underline{\underline{\Lambda}}_m|} \sum\limits_{i=1}^{N} \sum\limits_{m=1}^{M} w_m^{(i)} - \dfrac{1}{2} \sum\limits_{i=1}^{N} \sum\limits_{m=1}^{M} w_m^{(i)} (\underline{x} - \underline{\mu}_m)(\underline{x} - \underline{\mu}_m)^T$

$\implies \dfrac{\text{adj}(\underline{\underline{\Lambda}}_m)}{|\underline{\underline{\Lambda}}_m|} \sum\limits_{i=1}^{N} \sum\limits_{m=1}^{M} w_m^{(i)} = \sum\limits_{i=1}^{N} \sum\limits_{m=1}^{M} w_m^{(i)} (\underline{x} - \underline{\mu}_m)(\underline{x} - \underline{\mu}_m)^T$

$\underline{\underline{\Lambda}}_m^{-1} = \dfrac{\text{adj}(\underline{\underline{\Lambda}}_m)}{|\underline{\underline{\Lambda}}_m|} = \underline{\underline{\Sigma}}_m$

$\therefore \underline{\underline{\Sigma}}_m = \dfrac{\sum\limits_{i=1}^{N} \sum\limits_{m=1}^{M} w_m^{(i)} (\underline{x} - \underline{\mu}_m)(\underline{x} - \underline{\mu}_m)^T}{\sum\limits_{i=1}^{N} \sum\limits_{m=1}^{M} w_m^{(i)}}$

Q2. (b) Classification accuracy on training set: 79.515%
      Classification accuracy on testing set: 79.66%

(c)



(d) Accuracy on training dataset: 87.865%; Accuracy on testing dataset: 87.57%