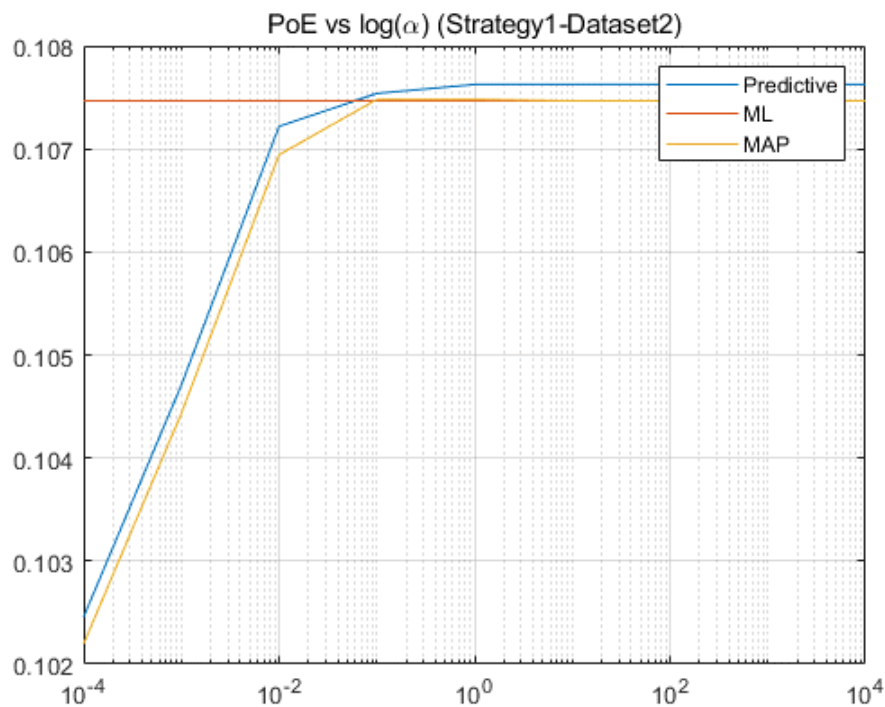
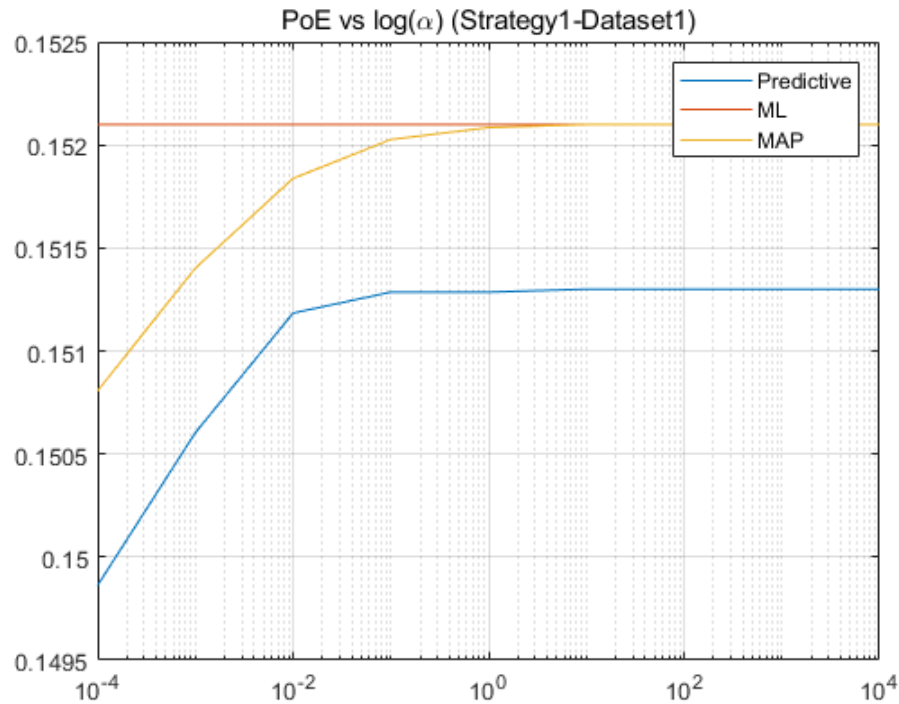


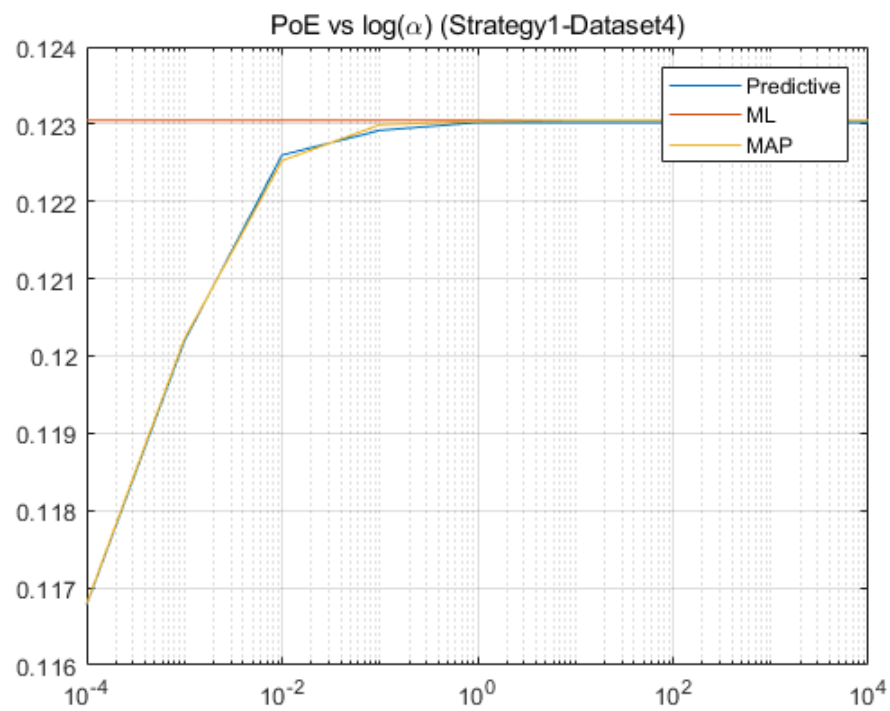
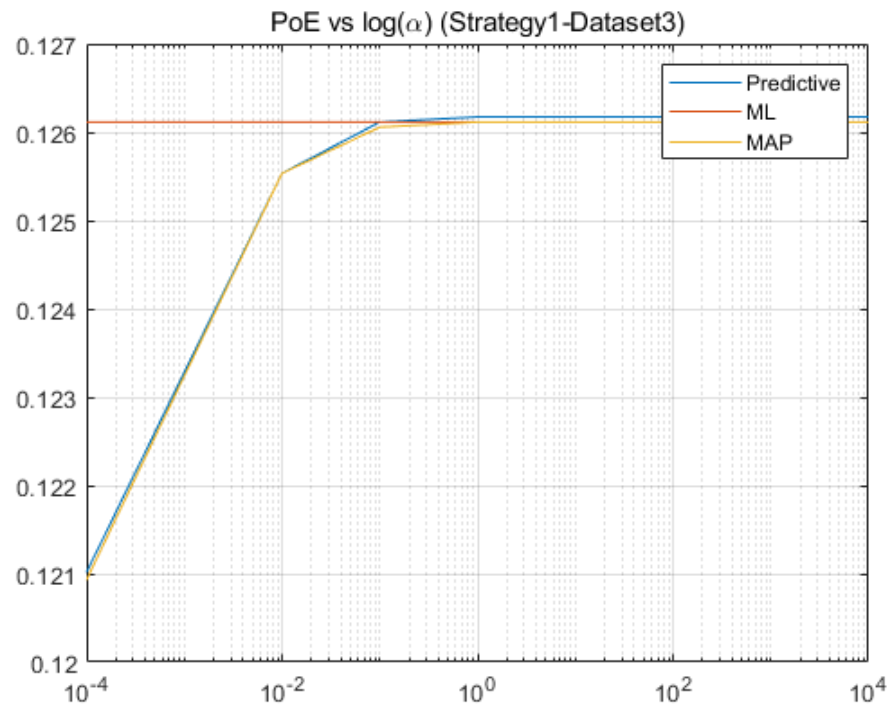
# [ECE271A] Homework 3 & 4 Solution

Name: Yang Yue PID: A53301503

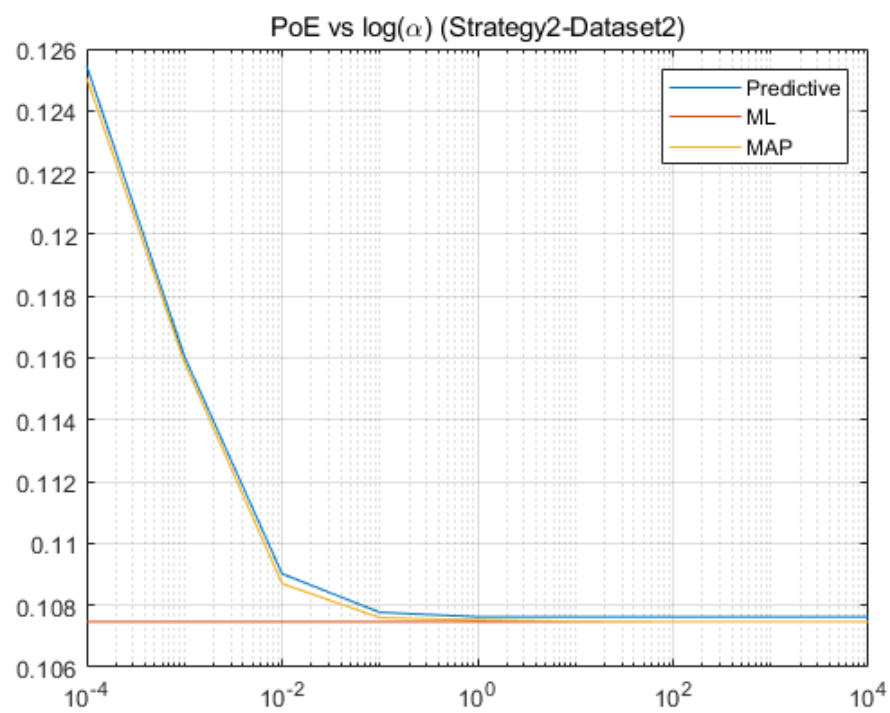
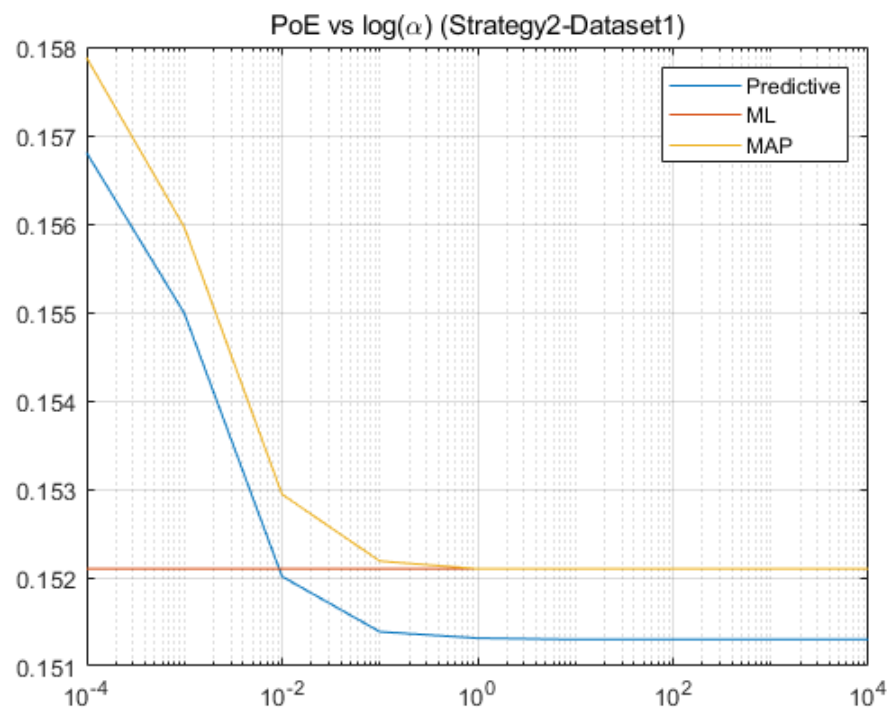
## Problem 4 (computer)

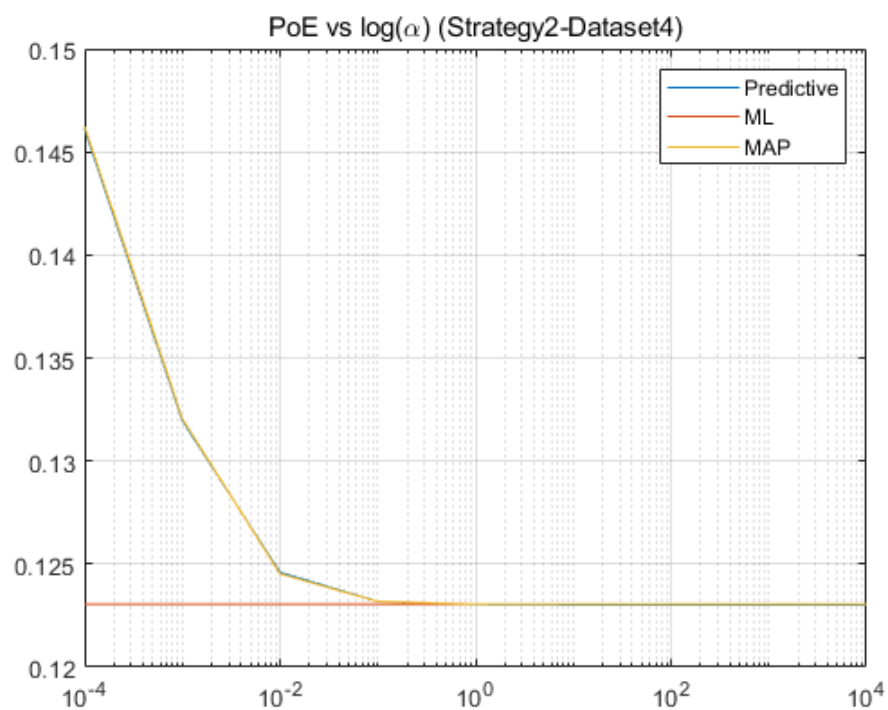
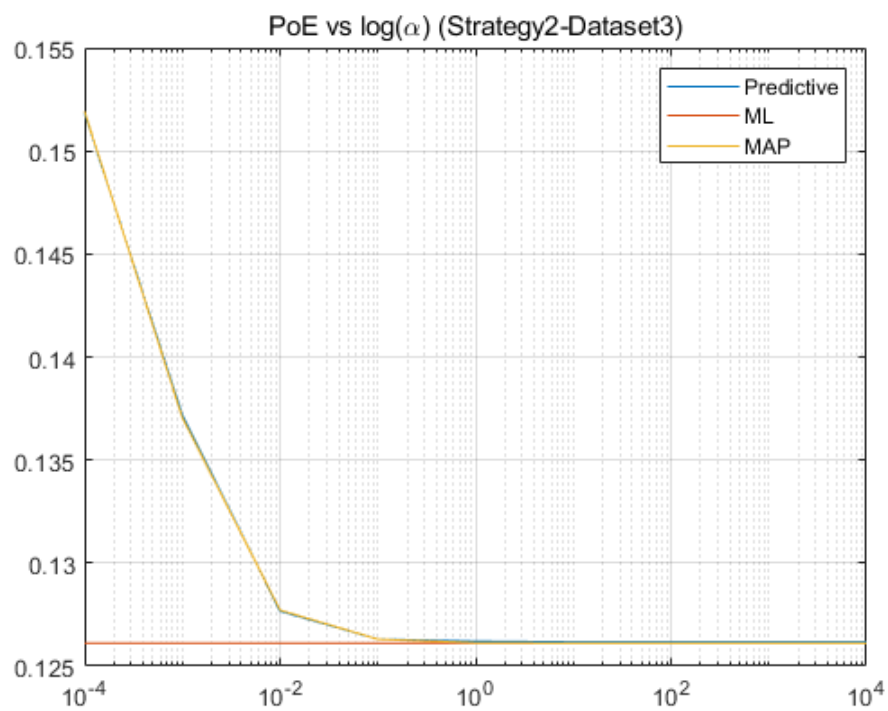
PoE (Probability of Error) vs different alphas with respect to different Datasets using Strategy 1





PoE (Probability of Error) vs different alphas with respect to different Datasets using Strategy 2





## Explanation

Given that

$$P_{\mu}(\mu) = \mathcal{G}(\mu, \mu_0, \Sigma_0),$$

where  $\mu_0$  is given,  $\Sigma_0^2 = \text{diag}(\alpha_i \cdot \omega_0)$  is given as well,

And

$$P_{\mu|T}(\mu|D) = \mathcal{G}(\mu, \mu_n, \Sigma_n),$$

$$P_{X|\mu, \Sigma} = \mathcal{G}(x, \mu, \Sigma)$$

Also, given that,

$$P_{X|T}(x|D) = \mathcal{G}(x, \mu_n, \Sigma + \Sigma_n)$$

In hw3, we assume that

$$\Sigma = \text{cov}(X), \text{ and } \Sigma_n = n\Sigma_0(\Sigma + n\Sigma_0)^{-1}.$$

For a given n

$$\mu_n = \underbrace{\frac{n\Sigma_0(\Sigma + n\Sigma_0)^{-1}}{\alpha_n}}_{\alpha_n} \mu_{ML} + \underbrace{\frac{\Sigma(\Sigma + n\Sigma_0)}{1-\alpha_n}}_{1-\alpha_n} \mu_0$$

Which can be written as,

$$\mu_n = \alpha_n \hat{\mu} + (1 - \alpha_0) \mu_0 \text{ where } \alpha_n \in [0, 1]$$

In MAP case,

$$P_{X|T}(x|D) = \mathcal{G}(x, \mu_{MAP}, \Sigma_n),$$

where

$$\mu_{MAP} = n\Sigma_0(\Sigma + n\Sigma_0)^{-1} \mu_{ML} + \Sigma(\Sigma + n\Sigma_0)^{-1} \mu_0,$$

In ML case,

$$P_{X|T}(x|D) = \mathcal{G}(x, \mu_{ML}, \Sigma_n),$$

- 1) In hw3,  $\sigma_0^2 = \alpha \times \omega_0$  represents how confident we are about the prior probability of  $\mu$ , for a hypothetical expectation of  $\mu$ , a larger  $\sigma_0^2$  shows less confidence.

Assuming that we have no confidence about the distribution of  $\mu$  is correct, we can show this by giving the distribution a very large covariance of distribution  $\Sigma_0$ . This indicates that  $\alpha_n \rightarrow 1$ , and  $\mu_n \rightarrow \hat{\mu} = \mu_{ML}$ , and because of that, the difference between of each curve using predictive equation, MAP, and MLE is tend to be very small when  $\alpha \rightarrow \infty$ .

Inversely, assuming that we are very certain about the distribution of  $\mu$ , we will make  $\Sigma_0$  be very small. Then  $\mu_n$  is determined by i) training dataset and ii)  $\mu_0$ .

As each plot of Strategy1 has shown, PoE increases with an increasing  $\alpha$ , which indicates that the assumption we make about the distribution of  $\mu$  is reasonable. Due to the reasonable assumption, intuitively, taking it into account can reduce the PoE. What the plots show is that a reasonable prior distribution of  $\mu$  help reduce the error rate, the more certainty we have about the correct assumption, the more "helpful" the assumption will be. It is obvious especially when the size of training set is small. As a result of this, PoE(predictive) is obviously smaller than PoE(MAP) and PoE(MLE) in Strategy1 with Dataset1.

- 2) When the size of training dataset is sufficiently large enough, let's suppose  $n \rightarrow \infty$ , then,

$$\Sigma_{Predictive} = \Sigma + \Sigma \Sigma_0 (\Sigma + n \Sigma_0)^{-1} = \Sigma + \frac{1}{n} \Sigma \Sigma_0 \left( \frac{1}{n} \Sigma + \Sigma_0 \right)^{-1} \rightarrow \Sigma = \Sigma_{MAP},$$

which implies that when  $\alpha$  remains, the larger the size of training set is, the closer the PoE(predictive) and PoE(MAP) get.

Also, while  $n \rightarrow \infty$ ,

$$\begin{aligned}\mu_{MAP} &= n\Sigma_0(n\Sigma_0^2 + \Sigma^2)^{-1}\mu_{ML} + \Sigma^2(n\Sigma_0^2 + \Sigma^2)\mu_0 \\ &= \Sigma_0\left(\Sigma_0^2 + \frac{1}{n}\Sigma^2\right)^{-1}\mu_{ML} + \frac{1}{n}\Sigma^2\left(\Sigma_0^2 + \frac{1}{n}\Sigma^2\right)\mu_0 \rightarrow \mu_{ML}\end{aligned}$$

Which indicates that as  $n$  gets larger, the plot of MAP solution and the plot of ML solution tend to coincide.

- 3) When the strategy changes from to Strategy2, it is suggesting that the distribution of  $P_\mu(\mu)$  of front ground and back ground are the same, thus unreasonable. Because we cannot distinguish  $\mu_{BG}$  and  $\mu_{FG}$ , which, with the intuition that there should be some differences of the distribution of parameter from 2 different things' images. It can be shown from the plots that as  $\alpha$  gets larger, the weight of prior of  $\mu$  gets less, the error rate decreases, which intuitively suggests that "the data is trying to fix the mistake made by wrong prior estimation".

Other relative behaviors remain even though the estimated distribution of prior is wrong. For example, i) the curves are still tending to be together with an increasing  $\alpha$ ; ii) as  $n$  grows, the estimation using Bayesian estimate combines the prior beliefs is more alike to the estimation using MAP rule.