

Method 3

Tuesday, November 9, 2021 12:55 PM

GIVEN: $P_{x|u, \hat{\Sigma}} = G(x, \mu, \hat{\Sigma})$ s.t. $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x_i - \frac{1}{N} \sum_{j=1}^N x_j)(x_i - \frac{1}{N} \sum_{j=1}^N x_j)^T$

$M = ?$

$P_u(u) = G(u, \mu_u, \hat{\Sigma}_u)$ s.t. strategy 1: $\begin{cases} \text{cheetah: } \mu_u = [1 \ 0 \ \dots \ 0]^T \\ \text{grass: } \mu_u = [0 \ 1 \ \dots \ 0]^T \\ \hat{\Sigma}_u = \text{diag}(d_{u1}) \text{ s.t. } u_i \text{ in Prior-3.mat} \end{cases}$

strategy 2: $\begin{cases} \text{cheetah and grass: } \mu_u = [0 \ 1 \ \dots \ 0]^T \\ \hat{\Sigma}_u = \text{diag}(d_{u1}) \text{ s.t. } u_i \text{ in Prior-2.mat} \end{cases}$

A) training set D_i : For each class, compute $\hat{\Sigma}$ formula above where x_i are samples in D_i

$\mu_i = \hat{\Sigma}_0 (\hat{\Sigma}_0 + \frac{1}{N} \hat{\Sigma})^{-1} \mu_{ML} + \frac{1}{N} \hat{\Sigma} (\hat{\Sigma}_0 + \frac{1}{N} \hat{\Sigma})^{-1} \mu_0$

$\hat{\Sigma}_1 = \hat{\Sigma}_0 (\hat{\Sigma}_0 + \frac{1}{N} \hat{\Sigma})^{-1} \hat{\Sigma}$

parameters of $P_{x|u}$ $\Rightarrow P_{x|u}(x|D_i) = G(x, \mu_i, \hat{\Sigma}_i)$

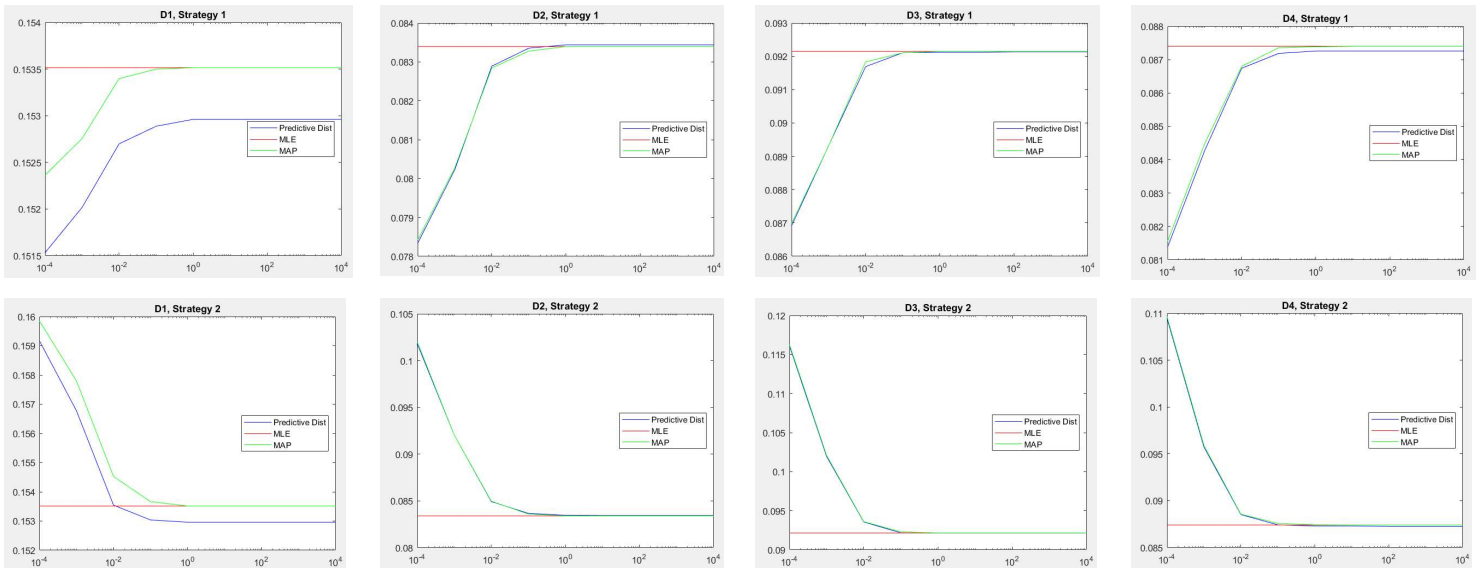
B) For ML procedure, we get: $\frac{P_{x|u}(x|D_i) \cdot G(x, \mu_{ML}, \hat{\Sigma}_{ML})}{\frac{1}{N} \sum_{i=1}^N x_i}$ $\hat{\Sigma}$ from formula above

C) $P_{x|u}(x|D_i) = P_{x|u}(x|\mu_{ML}) \Rightarrow \mu_{ML} = \arg\max_{\mu} P_{ML}(u|D_i)$

$= \arg\max_{\mu} G(u, \mu, \hat{\Sigma}_i)$

$= \mu_i$ by definition of Gaussian

$= G(x, \mu_i, \hat{\Sigma}_i)$



- MLE: μ_{ML} and $\hat{\Sigma}$ are not dependent on d , and hence the error is constant.
- MAP: μ_i is dependent on d ; When $d \rightarrow \infty$, the diagonal elements of $\hat{\Sigma}_0 \rightarrow \infty \Rightarrow \mu_i = \hat{\Sigma}_0 (\hat{\Sigma}_0 + \frac{1}{N} \hat{\Sigma})^{-1} \mu_{ML} + \frac{1}{N} \hat{\Sigma} (\hat{\Sigma}_0 + \frac{1}{N} \hat{\Sigma})^{-1} \mu_0 = \mu_{ML}$
- PD: $\mu_i \rightarrow \mu_{ML}$ as $d \rightarrow \infty$ as shown for MAP.

$\hat{\Sigma}_1 \rightarrow \frac{1}{N} \hat{\Sigma}$ as $d \rightarrow \infty$, so the variance for the predictive distribution: $\hat{\Sigma} \rightarrow (1 + \frac{1}{N}) \hat{\Sigma}$.

Hence, for large n , this distribution is equal to the MAP (D_1, D_2, D_3). For smaller $n(D_4)$, we have a higher variance.

- GENERAL: when $d \rightarrow \infty$, $\hat{\Sigma}_0 \rightarrow \infty$ meaning that the variance of the prior $\rightarrow \infty$ and hence the prior is not helpful in finding the solution, so we rely on the data. We rely on prior when $d \rightarrow 0$ and n is small

- in strategy 1, the prior is accurate meaning the more we rely on it ($d \rightarrow 0$), the better the error.
- in strategy 2, the prior is non-informative meaning the less we rely on it ($d \rightarrow \infty$), the better the error.
- D_1 generates the highest error (smallest n), and D_2, D_3, D_4 are all about the same magnitude

