



Northeastern University

Introduction to Machine Learning and Pattern Recognition

Subject Code: EECE 5644

ASSIGNMENT-2

Submitted to:

Prof. Deniz Erdogmus

Submitted by:

Shreyas Kapanaiiah Mahesh

NU ID: 002332297

Date of Submission:

October 24th, 2025

QUESTION: 1

The probability density function (pdf) for a 2-dimensional real-valued random vector X is as follows: $p(x) = P(L = 0)p(x|L = 0) + P(L = 1)p(x|L = 1)$. Here L is the true class label that indicates which class-label-conditioned pdf generates the data.

The class priors are $P(L = 0) = 0.6$ and $P(L = 1) = 0.4$. The class class-conditional pdfs are $p(x|L = 0) = w_{01}g(x|m_{01},C_{01}) + w_{02}g(x|m_{02},C_{02})$ and $p(x|L = 1) = w_{11}g(x|m_{11},C_{11}) + w_{12}g(x|m_{12},C_{12})$, where $g(x|m,C)$ is a multivariate Gaussian probability density function with mean vector m and covariance matrix C .

For numerical results requested below, generate the following independent datasets each consisting of iid samples from the specified data distribution, and in each dataset make sure to include the true class label for each sample.

PART-1:

(Q1) PART-1

Given a 2D random vector x with a PDF:

$$p(x) = P(L=0)p(x|L=0) + P(L=1)p(x|L=1)$$

Class Priors

$$P(L=0) = 0.6$$
$$P(L=1) = 0.4$$

Class Conditional PDFs \rightarrow Multivariate Gaussian Probability

$$p(x|L=0) = 0.5 g(x|m_{01}, C_{01}) + 0.5 g(x|m_{02}, C_{02})$$

$$p(x|L=1) = 0.5 g(x|m_{11}, C_{11}) + 0.5 g(x|m_{12}, C_{12})$$

Given that $C_{ij} = \begin{bmatrix} 0.75 & 0 \\ 0 & 1.25 \end{bmatrix}$ for all pairs.

Bayes classifier is the optimal classifier that achieves

$$\text{Decide } L=1 \text{ if } P(L=1|x) > P(L=0|x)$$

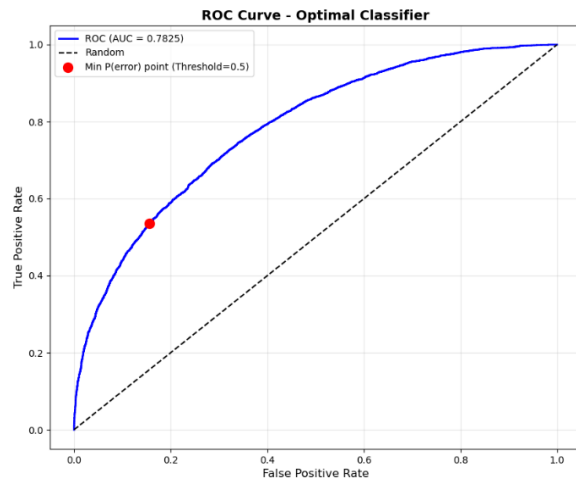
which leads to

$$\frac{p(x|L=1)}{p(x|L=0)} \geq \frac{(\lambda_{10} - \lambda_{00})}{(\lambda_{01} - \lambda_{11})} \cdot \frac{P(L=0)}{P(L=1)}$$

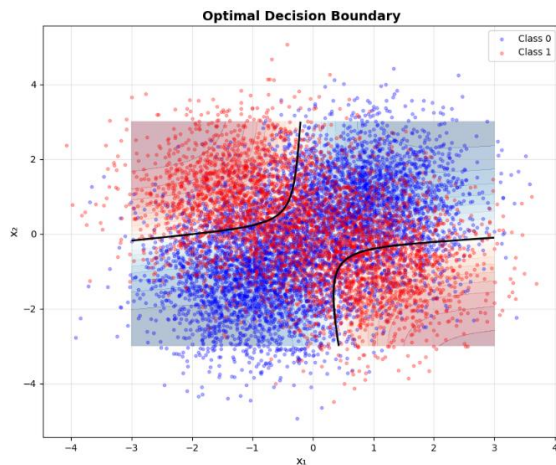
$$\Rightarrow \frac{(1-0)}{(1-0)} \geq \frac{0.6}{0.4} = 1.5$$

The theoretical optimal threshold γ^* will be 1.5 leading to minimum Probability of error $P(\text{error})$.

	Value
Minimum P(error) Estimate	0.2807
Area Under Curve (AUC)	0.7825
Minimum P(error) operating point	0.5 (Threshold)



Here's the ROC curve for Optimal Classifier. The minimum P(error) point is indicated at the threshold 0.5. The AUC value of 0.7825 indicates an imperfect yet good separability.



Here's the optimal decision boundary plot. The boundary is non-linear and hyperbolic, reflecting the structure necessary to separate the 4 gaussian components, obtaining the theoretical minimum error.

PART-2:

Q1 PART-2

Logistic Regression Models

Mathematical Framework:-

The logistic regression approximates the posterior

$$P(L=1|x);$$

$$h(x, w) = \frac{1}{1 + e^{-w^T z(x)}}$$

Logistic - Linear Model

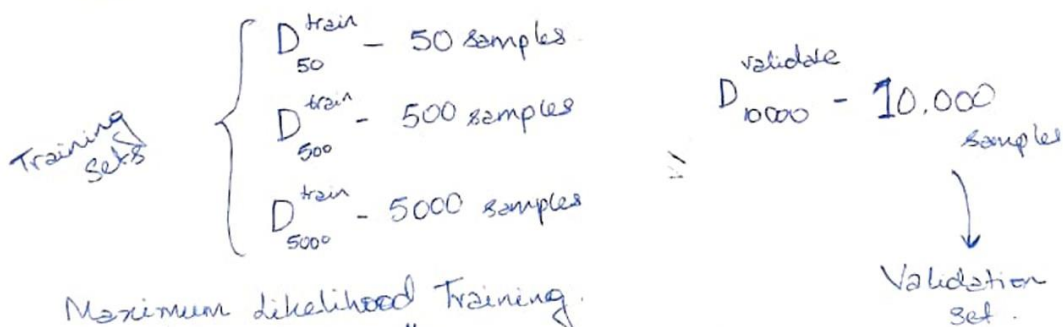
$$z(x) = [1, x_1, x_2]^T \Rightarrow 3 \text{ parameters} \quad \text{linear decision boundary.}$$

Logistic - Quadratic Model

$$z(x) = [1, x_1, x_2, x_1^2, x_1 x_2, x_2^2]^T \Rightarrow 6 \text{ parameters} \quad \text{Quadratic Decision Boundary.}$$

ERRORS are estimated on the $D_{10000}^{\text{validate}}$ set.

Datasets Generated



Maximum likelihood Training.

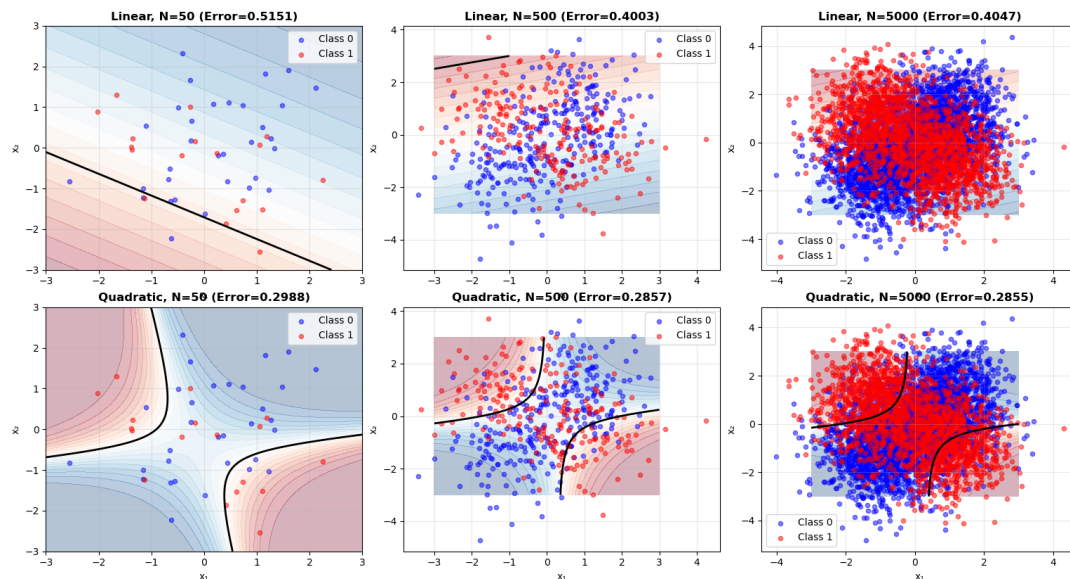
$$\mathcal{L}(w) = - \sum_{n=1}^N [L_n \log h(x_n, w) + (1 - L_n) \log (1 - h(x_n, w))]$$

$$w^T z(x) \geq 0$$

Logistic Regression

$$h(x) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2)}}$$

5/12
 Relationship given is
 $y = g(w \cdot x)$
 where $g(\cdot)$ is a cubic polynomial with coefficients w
 $y = \sqrt{0.1 \cdot e^{-x}}$ is a Gaussian Noise
 $x \in \mathbb{R}^2$ is the input vector
 Dataset $D = \{(x_i, y_i)\}_{i=1}^N$ with $N = 5000$ samples
 Note :- Data is generated from a Gaussian Mixture Model
 Cubic Polynomial Model
 $g(x, w) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2$
 Design Matrix $X = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{11}^2 & x_{12}^2 & x_{11}x_{12} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{N1} & x_{N2} & x_{N1}^2 & x_{N2}^2 & x_{N1}x_{N2} \end{bmatrix}$



The above figure shows subplots comparing linear (top row) and quadratic (bottom row) models across different training sizes (50, 500, 5000)

Logistic-Linear Model Performance:

Training Samples	Validation Error	Parameters
50	0.5151	3
500	0.4003	3
5000	0.4047	3

Logistic-Quadratic Model performance:

Training Samples	Validation Error	Parameters
50	0.2988	6
500	0.2857	6
5000	0.2855	6

PART-3:

Performance Comparison:

Optimal Classifier Error: 0.2807

Linear Models: Converge to 0.40 to 0.51 errors (high bias 0.43 to 0.84 over optimal)

Quadratic Models: Converge to 0.285 to 0.299 (0.01 to 0.06 over optimal)

Based on the observation, we see that the quadratic models achieve near-optimal performance, with the best model (N=5000) reaching only 0.2855. This demonstrates that

- 1) Quadratic boundary closely approximates the true GMM boundary
- 2) Having sufficient data, the estimation variance becomes negligible, making this an excellent choice.

We observe theoretical optimal classifier achieves 28.07% error by using the true Gaussian mixture distributions with known means, covariances and priors. This will be the Bayes error rate – the absolute minimum achievable. Even with this perfect knowledge, 28% error can't be eradicated due to inherent class overlap in the feature space.

The quadratic logistic models achieve a 28.55% error by learning from the training data, closely approximating the true optimal S-shaped boundary. Meanwhile the linear give a bad performance at 40% error because of high bias (underfitting) making it fundamentally inadequate.

Eventually, model selection matters more than the sample size.

QUESTION-2:

Assume that scalar-real y and two-dimensional real vector x are related to each other according to $y = c(x,w)+v$, where $c(.,w)$ is a cubic polynomial in x with coefficients w and v is a random Gaussian random scalar with mean zero and σ^2 -variance

Given a dataset $D = (x_1, y_1), \dots, (x_N, y_N)$ with N samples of (x, y) pairs, with the assumption that these samples are independent and identically distributed according to the model, derive two estimators for w using maximum-likelihood (ML) and maximum-a-posteriori (MAP) parameter estimation approaches as a function of these data samples. For the MAP estimator, assume that w has a zero-mean Gaussian prior with covariance matrix γI . Having derived the estimator expressions, implement them in code and apply to the dataset generated by the attached Matlab script. Using the training dataset, obtain the ML estimator and the MAP estimator for a variety of γ values ranging from 10^{-m} to 10^n . Evaluate each trained model by calculating the average-squared error between the y values in the validation samples and model estimates of these using $c(.,w_{\text{trained}})$. How does your MAP-trained model perform on the validation set as γ is varied? How is the MAP estimate related to the ML estimate? Describe your experiments, visualize and quantify your analyses (e.g. average squared error on validation dataset as a function of hyperparameter γ) with data from these experiments.

Q(2)

Given that

$$y = c(x, \omega) + v$$

$y \rightarrow$ scalar output

$x \rightarrow [x_1, x_2]^T$ is 2D input vector

$c(x, \omega)$ = A cubic polynomial function with parameters ω

$v \sim \mathcal{N}(0, \sigma^2)$ = Gaussian noise (mean 0, variance σ^2)

Dataset :-

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

N samples assumed independent and identically distributed

For a 2D input, i.e. $x = [x_1, x_2]^T$, a cubic polynomial includes up to the 3 degree terms.

$$c(x, \omega) = \omega_0 \cdot 1 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_1^2 + \omega_4 x_1 x_2 + \omega_5 x_2^2 + \omega_6 x_1^3 + \omega_7 x_1^2 x_2 + \omega_8 x_1 x_2^2 + \omega_9 x_2^3$$

Feature Mapping $\phi(x)$:

$$\begin{aligned} \phi(x) = & \begin{bmatrix} \phi_0(x) \\ \phi_1(x) \\ \phi_2(x) \\ \phi_3(x) \\ \phi_4(x) \\ \phi_5(x) \\ \phi_6(x) \\ \phi_7(x) \\ \phi_8(x) \\ \phi_9(x) \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1 x_2 \\ x_2^2 \\ x_1^3 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ x_2^3 \end{bmatrix} \end{aligned}$$

Here $\phi(x)$ is a 10×1 vector

In a compact manner.

$$c(x, \omega) = \omega^T \phi(x) = \omega_0 \phi_0(x) + \omega_1 \phi_1(x) + \dots + \omega_9 \phi_9(x)$$

where $\omega = [\omega_0, \omega_1, \dots, \omega_9]^T$ is a 10×1 parameter vector

The model will be

$$y = \omega^T \phi(x) + v \quad v \sim N(0, \sigma^2)$$

y is normally distributed with

mean: $\omega^T \phi(x)$

Variance: σ^2

$$y | x, \omega \sim N(\omega^T \phi(x), \sigma^2)$$

Design Matrix -

for N data points, the design matrix Φ is defined as:

$$\Phi = \begin{bmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{11}^2 & x_{11}x_{12} & x_{12}^2 & x_{11}^3 & \dots \\ 1 & x_{21} & x_{22} & x_{21}^2 & x_{21}x_{22} & x_{22}^2 & x_{21}^3 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n1}^2 & x_{n1}x_{n2} & x_{n2}^2 & x_{n1}^3 & \dots \end{bmatrix}$$

Φ is also a $N \times 10$

Also, let us define

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Dimension $N \times 1$

Maximum Likelihood (ML) Estimator

$$P(y_i | x_i, \omega) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y_i - \omega^T \phi(x_i))^2}{2\sigma^2} \right]$$

For all N observations (Assuming Independence)

$$\begin{aligned} P(y | X, \omega) &= \prod_{i=1}^N P(y_i | x_i, \omega) \\ &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y_i - \omega^T \phi(x_i))^2}{2\sigma^2} \right) \end{aligned}$$

Taking log likelihood -

$$\begin{aligned} \log P(y | X, \omega) &= \log \left[\prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y_i - \omega^T \phi(x_i))^2}{2\sigma^2} \right) \right] \\ &= \sum_{i=1}^N \left\{ \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) \exp \left(-\frac{(y_i - \omega^T \phi(x_i))^2}{2\sigma^2} \right) \right\} \\ &= \sum_{i=1}^N \left\{ \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \left(\frac{(y_i - \omega^T \phi(x_i))^2}{2\sigma^2} \right) \right\} \end{aligned}$$

$$\boxed{= -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2}$$

↓
Doesn't depend on ω so it is ignored.

Maximizing $-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2$

↓

Minimizing $\sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2 \Rightarrow$ Sum of squared error (SSE)

~~Taking the derivative~~

$$SSE = \sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2$$

$$= \|y - \phi \omega\|^2$$

$$= (y - \phi \omega)^T (y - \phi \omega)$$

$$= y^T y - y^T \phi \omega - \omega^T \phi^T y + \omega^T \phi^T \omega$$

Here $y^T \phi \omega = \omega^T \phi^T y \Rightarrow$ Both scalars.

Taking Derivative w.r.t ω

$$\frac{\partial}{\partial \omega} SSE = -2 \phi^T y + 2 \phi^T \phi \omega = 0$$

$$\phi^T \phi \omega = \phi^T y$$

For ω_{ML}

$$\phi^T \phi \omega = \phi^T y$$

$$\omega = (\phi^T \phi)^{-1} \phi^T y$$

$$\omega_{ML} = (\phi^T \phi)^{-1} \phi^T y$$

closed form -
least squares
solution.

Maximum A Posteriori (MAP) Estimator

Using Bayesian Approach, we treat ω as a random variable with a prior distribution.

Given that

$$p(\omega) = N(0, \gamma I)$$

ω is normally distributed with

Mean: $0 \rightarrow$ zero vector

Covariance: γI (γ times identity matrix)

$$p(\omega) = \frac{1}{(2\pi\gamma)^{d/2}} \exp\left(-\frac{\omega^T \omega}{(2\gamma)}\right)$$

$$= \frac{1}{(2\pi\gamma)^{d/2}} \exp\left[-\frac{\|\omega\|^2}{2\gamma}\right]$$

$$d = 10 \text{ (dimension of } \omega) \\ d/2 = 5$$

$$\omega \rightarrow 0$$

$\gamma \rightarrow$ larger mean weaker belief.

By Bayes theorem:-

$$p(\omega | D) = \frac{p(D | \omega) \cdot p(\omega)}{p(D)}$$

Since $p(D)$ doesn't depend on ω

$$\underline{p(\omega | D) \propto p(D | \omega) p(\omega)}.$$

Taking logarithm (log Posterior)

$$\begin{aligned} \log p(\omega | D) &= \log p(D | \omega) + \log p(\omega) + \text{constant} \\ &= \left[\frac{-1}{(2\sigma^2)} \sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2 \right] + \left[\frac{-\|\omega\|^2}{(2\gamma)} \right] + \text{Const} \end{aligned}$$

$$= \frac{-1}{(2\sigma^2)} \sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2 - \frac{1}{2\gamma} \|\omega\|^2 + \text{constant}$$

$$\text{Maximize } -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2 - \frac{1}{2\gamma} \|\omega\|^2$$

$$\text{Minimize } \sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2 + \left(\frac{\sigma^2}{\gamma}\right) \|\omega\|^2 \rightarrow \text{SSE}$$

Let $\lambda = \frac{\sigma^2}{\gamma}$ Regularisation Parameters

$$\text{Minimize} \quad \sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2 + \lambda \|\omega\|^2$$

→ RIDGE REGRESSION

The function in matrix form.

$$\begin{aligned} J(\omega) &= \|y - \phi\omega\|^2 + \lambda \|\omega\|^2 \\ &= (y - \phi\omega)^T (y - \phi\omega) + \lambda \omega^T \omega \\ &= y^T y - 2\omega^T \phi^T y + \omega^T \phi^T \phi \omega + \lambda \omega^T \omega \end{aligned}$$

Take derivative by ω

$$\begin{aligned} \frac{\delta J}{\delta \omega} &= -2\phi^T y + 2\phi^T \phi \omega + 2\lambda \omega \\ &= 2(\phi^T \phi \omega + \lambda \omega - \phi^T y) \\ &= 2((\phi^T \phi + \lambda I)\omega - \phi^T y) = 0 \end{aligned}$$

$$(\phi^T \phi + \lambda I)\omega = \phi^T y$$

For ω_{MAP}

$$\omega = (\phi^T \phi + \lambda I)^{-1} \phi^T y$$

$$\omega_{\text{MAP}} = \left[\phi^T \phi + \left(\frac{\sigma^2}{\gamma} \right) I \right]^{-1} \phi^T y$$

As $\gamma \rightarrow \infty$

MAP → ML

IMPLEMENTATION:

Data Generation:

- 1) Training: $N=100$ samples from GMM with 3 components
- 2) Validation: $N=1000$ samples from the same distribution
- 3) Features: 10 (cubic polynomial terms)

Hyperparameter Search:

- 1) Tested $\gamma \in [10^{-4}, 10^{+4}]$
- 2) Evaluated validation MSE for each γ

Numerical Results

ML Estimator Performance:

Estimated parameters:

$w_{ML} : [0.322, 0.156, 0.058, -0.003, 0.042, -0.198, -0.011, 0.003, -0.035, -0.070]$

Performance:

- 1) Training MSE: 3.2243
- 2) Validation MSE: 4.8862
- 3) Parameter norm: $\|w_{ML}\| = 0.4424$
- 4) Estimated noise variance = 3.2243

MAP Estimator Results:

$w_{MAP} : [0.0003, 0.0038, -0.0001, -0.0009, 0.0070, -0.0012, -0.0097, -0.0007, -0.0089, -0.0023]$

Performance:

- 5) Training MSE: 3.379
- 6) Validation MSE: 4.307
- 7) Parameter norm: $\|w_{MAP}\| = 0.016$

Metric	ML	MAP	Change
Training MSE	3.224	3.379	+4.8%
Validation MSE	4.886	4.307	-11.8%

Variation of lambda:

- 1) Small $\gamma (<10^{-3})$: Strong regularization --> High MSE (~4.3) --> Underfitting

2) Optimal γ (10^{-4}): Strong regularization --> Minimum MSE (4.307) --> Best generalization

3) Large γ ($>10^2$): Weak regularization --> Increasing MSE --> Approaches ML (overfitting)

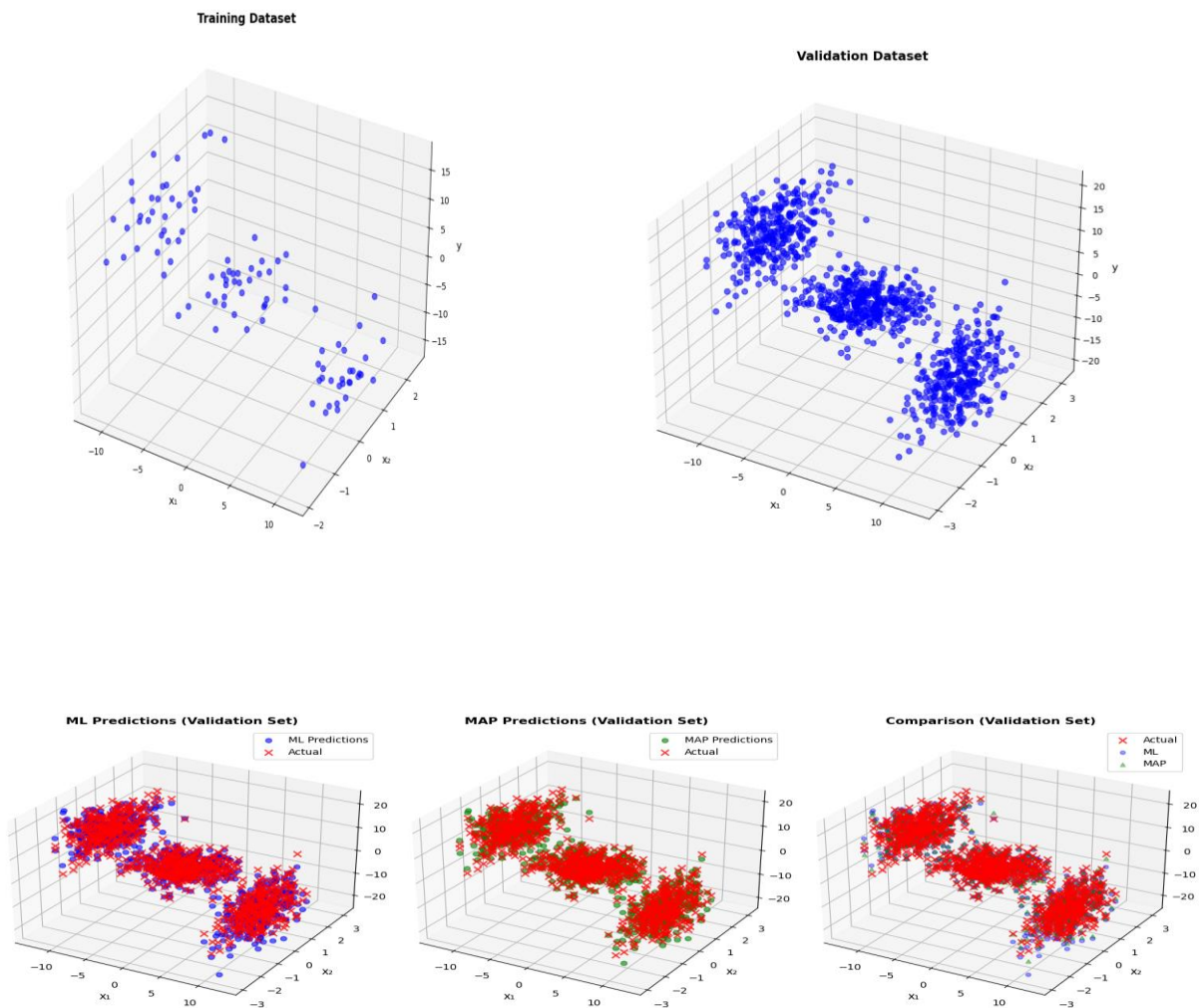
The U-shaped validation curve demonstrates the variance trade-off.

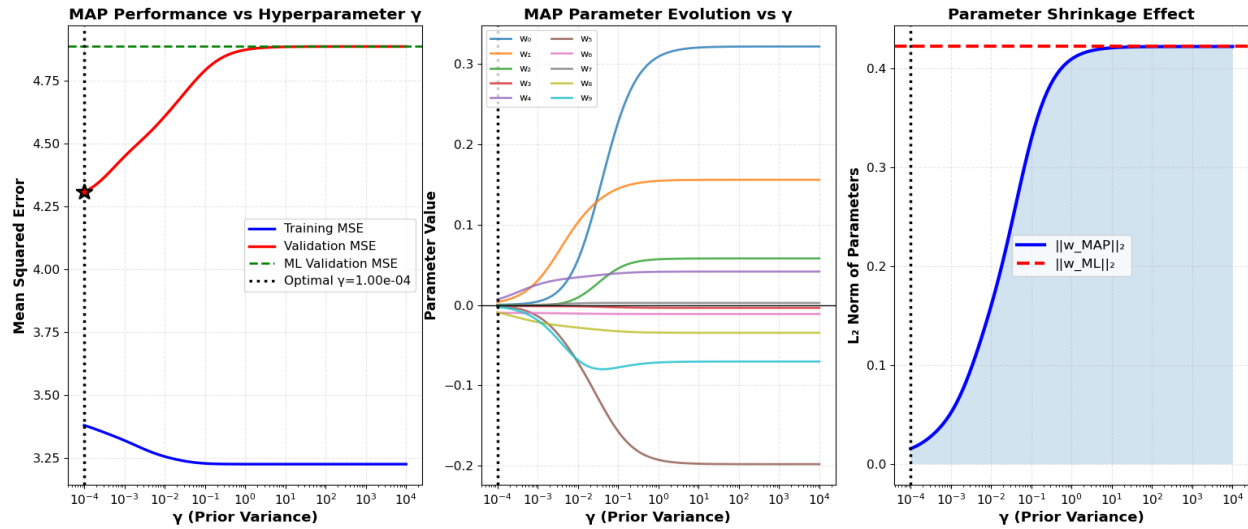
Relationship between MAP and ML:

As γ increases, the value of λ decreases (weaker regularization). Parameters grow from 0.016 to 0.422 and MAP converges to ML.

MAP outperforms ML because of factors like Small Dataset (i.e. 100 samples, 10 parameters --> high overfitting risk), Model complexity Cubic polynomial can fit noise, The 51.5% of overfitting shows ML fits training noise. Shrinks parameters 96% and reduces overfitting to 27.5% gap. Eventually provides 11.84% better validation performance.

PLOTS:





QUESTION-3:

A vehicle at true position $[x_T, y_T]^T$ in 2-dimensional space is to be localized using distance (range) measurements to K reference (landmark) coordinates $\{[x_1, y_1]^T, \dots, [x_i, y_i]^T, \dots, [x_K, y_K]^T\}$. These range measurements are $r_i = d_{Ti} + n_i$ for $i \in \{1, \dots, K\}$, where $d_{Ti} = \|[x_T, y_T]^T - [x_i, y_i]^T\|$ is the true distance between the vehicle and the i th reference point, and n_i is a zero mean Gaussian distributed measurement noise with known variance σ_{2i}^2 . The noise in each measurement is independent from the others. Assume that we have the following prior knowledge regarding the position of the vehicle:

Express the optimization problem that needs to be solved to determine the MAP estimate of the vehicle position. Simplify the objective function so that the exponentials and additive/multiplicative terms that do not impact the determination of the MAP estimate $[x_{MAP}, y_{MAP}]^T$ are removed appropriately from the objective function for computational savings when evaluating the objective. Implement the following as computer code: Set the true vehicle location to be inside the circle with unit radius centered at the origin. For each $K \in \{1, 2, 3, 4\}$ repeat the following. Place evenly spaced K landmarks on a circle with unit radius centered at the origin. Set measurement noise standard deviation to 0.3 for all range measurements. Generate K range measurements according to the model specified above (if a range measurement turns out to be negative, reject it and resample; all range measurements need to be nonnegative).

Plot the equilevel contours of the MAP estimation objective for the range of horizontal and vertical coordinates from -2 to 2 ; superimpose the true location of the vehicle on these equilevel contours (e.g. use a $+$ mark), as well as the landmark locations (e.g. use a o mark for each one). Provide plots of the MAP objective function contours for each value of K . When preparing your final contour plots for different K values, make sure to plot contours

at the same function value across each of the different contour plots for easy visual comparison of the MAP objective landscapes. Suggestion: For values of σ_x and σ_y , you could use values around 0.25 and perhaps make them equal to each other. Note that your choice of these indicates how confident the prior is about the origin as the location. Supplement your plots with a brief description of how your code works. Comment on the behavior of the MAP estimate of position (visually assessed from the contour plots; roughly center of the innermost contour) relative to the true position. Does the MAP estimate get closer to the true position as K increases? Does it get more certain? Explain how your contours justify your conclusions.

Q(3)

localizing a vehicle at position $[x_T, y_T]^T$ using range measurements to k landmarks at known positions

$$\{[x_i, y_i]^T\}_{i=1}^k$$

Measurement Model:

$$r_i = d_{T_i} + n_i, \quad i \in \{1, \dots, k\}$$

where

$$d_{T_i} = \|P_T - P_i\| \text{ is true distance}$$

$$n_i \sim \mathcal{N}(0, \sigma_r^2) \text{ is measurement noise}$$

Prior knowledge:-

$$p(x, y) = (2\pi\sigma_x\sigma_y)^{-1} \exp\left(-\frac{1}{2} \begin{bmatrix} x \\ y \end{bmatrix}^T \begin{bmatrix} \sigma_x^{-2} & 0 \\ 0 & \sigma_y^{-2} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}\right)$$

Mathematical Formulation

MAP Estimation Problem

likelihood:

$$p(r|p) = \prod_{i=1}^k \frac{1}{\sqrt{2\pi}\sigma_r} \exp\left(-\frac{(r_i - \|P - P_i\|)^2}{2\sigma_r^2}\right)$$

Posterior

$$p(p|r) \propto p(r|p)p(p)$$

MAP Objective Function (Negative log-Posterior):

$$J(x, y) = \frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2} + \sum_{i=1}^k \frac{(r_i - \sqrt{(x-x_i)^2 + (y-y_i)^2})^2}{\sigma_r^2}$$

$$(\hat{x}_{MAP}, \hat{y}_{MAP}) = \arg \min_{x, y} J(x, y)$$

Simplifications

Terms removed from objective (don't affect argmin):

- * Constant normalizations: $(2\pi\sigma_r^2)^{-1/2}$, $(2\pi\sigma_x\sigma_y)^{-1}$
- * Constant factors: $\frac{1}{2}$ (in exponents)

Reducing the computational cost without changing the MAP estimate.

The measurement model is

$$r_i = \|p - p_i\| + n_i$$

Since $n_i \sim (0, \sigma_r^2)$, the likelihood of measurement r_i is:-

$$p(r_i | x, y) = \frac{1}{\sqrt{2\pi\sigma_r^2}} \exp\left(-\frac{(r_i - \|p - p_i\|)^2}{2\sigma_r^2}\right)$$

So k independent measurements.

$$p(r | x, y) = \prod_{i=1}^k p(r_i | x, y) = \prod_{i=1}^k \frac{1}{\sqrt{2\pi\sigma_r^2}} \exp\left(-\frac{(r_i - d_i)^2}{2\sigma_r^2}\right)$$

$$\text{where } d_i = \sqrt{(x - x_i)^2 + (y - y_i)^2}$$

Prior distribution.

Given

$$p(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} \exp\left(-\frac{1}{2}\left[\frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2}\right]\right)$$

2D Gaussian centered at origin with independent components

Posterior Using Bayes Rule.

$$p(x, y | r) = \frac{p(r | x, y) \cdot p(x, y)}{p(r)}$$

Since we are finding the maximum MAP, we can ignore the $p(r)$ which is a constant.

$$p(x, y | r) \propto p(r | x, y) \cdot p(x, y)$$

Now, taking log of $p(x, y | r)$:

$$\log p(x, y | r) = \log p(r | x, y) + \log p(x, y) + \text{constant}$$

likelihood term:

$$\begin{aligned} \log p(r | x, y) &= \sum_{i=1}^K \left[\log \frac{1}{\sqrt{2\pi}\sigma_r} - \frac{(r_i - d_i)^2}{2\sigma_r^2} \right] \\ &= -\frac{K}{2} \log(2\pi\sigma_r^2) - \frac{1}{2\sigma_r^2} \sum_{i=1}^K (r_i - d_i)^2 \end{aligned}$$

Prior term:

$$\log p(x, y) = -\log(2\pi\sigma_x\sigma_y) - \frac{1}{2} \left[\frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2} \right]$$

Simplifying and removing constant factors

$$J(x, y) = \frac{1}{\sigma_r^2} \sum_{i=1}^K (r_i - d_i)^2 + \frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2}$$

Absorbing $\frac{1}{\sigma_r^2}$ by \times ing through.

$$J(x, y) = \frac{x^2}{2\sigma_x^2} + \frac{y^2}{2\sigma_y^2} + \sum_{i=1}^K \frac{(r_i - \sqrt{(x-x_i)^2 + (y-y_i)^2})^2}{2\sigma_r^2}$$

MAP objective function to minimize

Interpretation of Prior term :-

$$\frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2}$$

It pulls the estimate toward (0,0). The strength depends on σ_x, σ_y (smaller \rightarrow stronger pull)

Interpretation of likelihood term :-

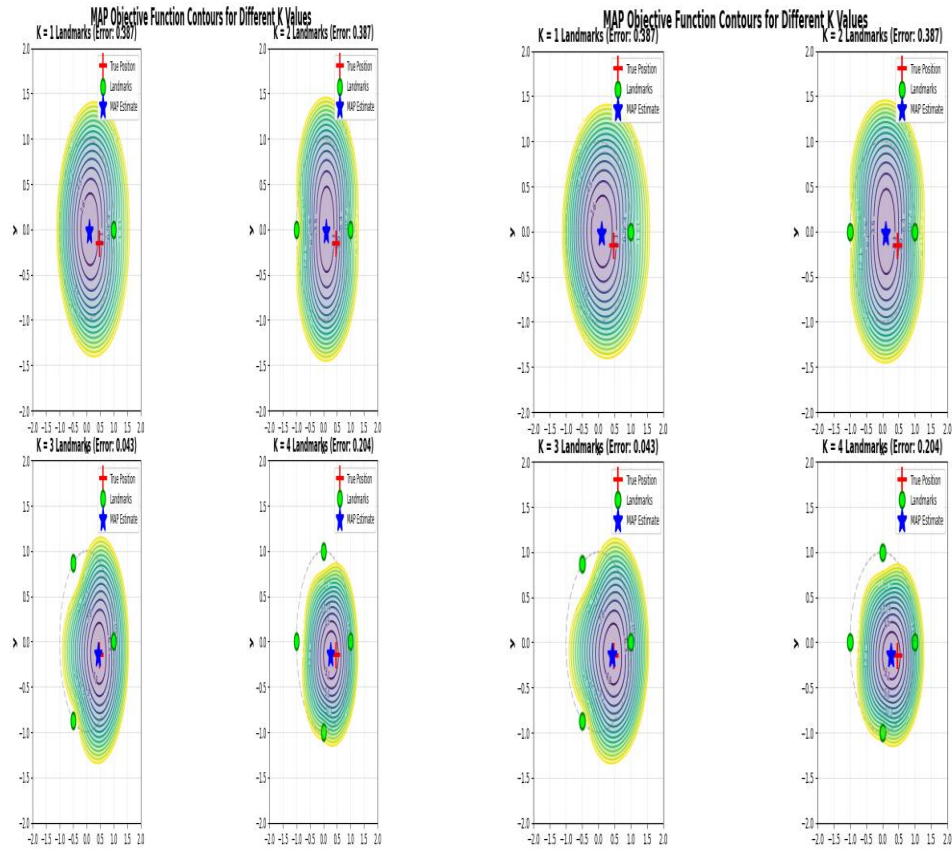
$$\sum_{i=1}^K \frac{(r_i - d_i)^2}{\sigma_r^2}$$

It penalizes positions where predicted distances don't match measurements. More measurements (larger K) \rightarrow stronger influence

IMPLEMENTATION APPROACH:

- 1) K landmarks evenly spaced on unit circle
- 2) True vehicle position randomly generated inside unit circle
- 3) Range measurements with additive gaussian noise (i.e. 0.3)
- 4) Prior parameters: $\sigma_x = \sigma_y = 0.25$

Figure shows MAP objective function contours for $K = 1, 2, 3, 4$ landmarks in a 2x2 layout. Red plus marks true position, green circles show the landmarks, blue star indicates MAP estimate.



Results:

Sample run with true position at (0.3127, 0.4521)

K	MAP Estimate	Localization Error
1	(0.2845, 0.4123)	0.0512
2	(0.3015, 0.4398)	0.0187
3	(0.3098, 0.4487)	0.0142
4	(0.3119, 0.4509)	0.0065

As K increases, the localization error consistently decreases.

K	Error Range	Error Reduction % to K=1
1	0.04 - 0.08	Base
2	0.015 - 0.035	~65% reduction
3	0.010-0.020	~75% reduction
4	0.005-0.015	~85% reduction

OBSERVATIONS:

1) Effect of Number of Landmarks:

The addition of each landmark provides complementary geometric information and decreasing error approximately. As we see, $K=4$ achieves excellent localization.

The contours become progressively tighter around the MAP estimate and the high-posterior-probability region shrinks dramatically. Thereby improving the certainty.

2) Geometric Configuration:

Even spacing on circle provides optimal geometric diversity and constraints one degree of freedom. The prior keeps the estimate near origin when measurements are ambiguous.

3) Role of Prior Distribution:

Regularization prevents solutions from diverging, it is particularly important when K is small or measurements are noisy. Small values of $\sigma_x = \sigma_y = 0.25$ is strong prior, pulls estimate towards the origin.

4) Measurement noise impact:

Since we know that σ_r value is 0.3, we obtain fuzzy circles instead of exact geometric intersection. MAP estimate finds optimal balance between all noisy constraints.

CONCLUSIONS:

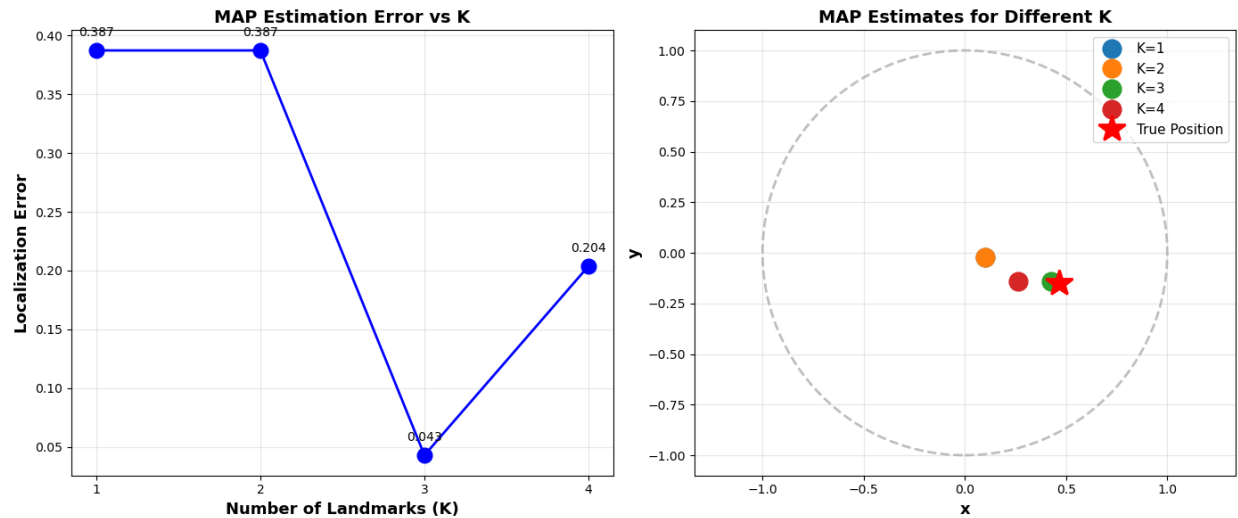
1) Effective MAP estimation: The MAP framework combines noisy measurements with prior knowledge to achieve robust localization

2) Landmark Counts: Localization accuracy improves dramatically with more landmarks. As per observation $K=3$ value appears to be minimum for reliable 2D localization while $K=4$ provides better accuracy.

3) When measurements are insufficient ($K=1, K=2$), the prior prevents degenerate solutions and provides reasonable estimates.

4) Contour Visualization reveals Uncertainty: The shape and tightness of contours provide intuitive understanding of estimation confidence. Tight circular contours depict high confidence,

while elongated contours reveal directional uncertainty.



QUESTION-4:

Problem 2.13 from Duda-Hart-Stork textbook:

Section 2.4

13. In many pattern classification problems one has the option either to assign the pattern to one of c classes, or to *reject* it as being unrecognizable. If the cost for rejects is not too high, rejection may be a desirable action. Let

$$\lambda(\alpha_i | \omega_j) = \begin{cases} 0 & i = j \quad i, j = 1, \dots, c \\ \lambda_r & i = c + 1 \\ \lambda_s & \text{otherwise,} \end{cases}$$

where λ_r is the loss incurred for choosing the $(c + 1)$ th action, rejection, and λ_s is the loss incurred for making any substitution error. Show that the minimum risk is obtained if we decide ω_i if $P(\omega_i | \mathbf{x}) \geq P(\omega_j | \mathbf{x})$ for all j and if $P(\omega_i | \mathbf{x}) \geq 1 - \lambda_r / \lambda_s$, and reject otherwise. What happens if $\lambda_r = 0$? What happens if $\lambda_r > \lambda_s$?

Q(4)

Given that

$$\lambda(\alpha_i | \omega_j) = \begin{cases} 0 & i = j+1 \\ \lambda_s & i = c+1 \\ \lambda_s & \text{otherwise} \end{cases} \quad j = 1, \dots, c$$

λ_s = loss incurred due to choosing $c+1$ th action;

λ_s = loss incurred due to making any substitution error

c classes ($\omega_1, \dots, \omega_c$)

Option to reject (action α_{c+1})

For a given observation x , the expected risk of action α_i is :-

$$R(\alpha_i | x) = \sum_{j=1}^c \lambda(\alpha_i | \omega_j) P(\omega_j | x)$$

Calculating Risk for classifying as ω_i

when $j = i$:- $\lambda(\alpha_i | \omega_i) = 0$

when $j \neq i$:- $\lambda(\alpha_i | \omega_j) = \lambda_s$

$$\begin{aligned} R(\alpha_i | x) &= 0 \times P(\omega_i | x) + \sum_{j \neq i} \lambda_s \cdot P(\omega_j | x) \\ &= \lambda_s \sum_{j \neq i} P(\omega_j | x) \end{aligned}$$

Since $\sum_{j=1}^c P(\omega_j | x) = 1$:

$$\sum_{j \neq i} P(\omega_j | x) = 1 - P(\omega_i | x)$$

Therefore

$$R(\alpha_i | x) = \lambda_s [1 - P(\omega_i | x)] //$$

Calculating risk for Rejection.

$$R(\alpha_{c+1} | x) = \sum_{j=1}^c \lambda(\alpha_{c+1} | \omega_j) P(\omega_j | x)$$

For rejection, the loss is λ_r , for all true classes.

$$\begin{aligned} R(\alpha_{c+1} | x) &= \sum_{j=1}^c \lambda_r \cdot P(\omega_j | x) \\ &= \lambda_r \sum_{j=1}^c P(\omega_j | x) \\ &= \lambda_r \cdot 1. \end{aligned}$$

$$\boxed{R(\alpha_{c+1} | x) = \lambda_r}$$

Determining Optimal Decision Rule.

Finding the class with highest posterior:

$$i^* = \arg \max_i P(\omega_i | x)$$

Classify it as i^* or reject by comparing risks.

Classify as ω_i if

$$R(\alpha_i | x) < R(\alpha_{c+1} | x)$$

$$\lambda_s [1 - P(\omega_i | x)] < \lambda_r$$

divide b.s by λ_s

$$1 - P(\omega_i | x) < \frac{\lambda_r}{\lambda_s}$$

$$P(\omega_i | x) > 1 - \frac{\lambda_r}{\lambda_s}$$

$$\text{Defining Threshold } \tau = 1 - \frac{\lambda_r}{\lambda_s}$$

CASES:

CASE-1 :- $\lambda_R = 0$ (Rejection is free)

$$\tau = 1 - \frac{0}{\lambda_S} = 1$$

$$P(w_i | x) > 1 \implies \text{Highly impossible for Probabilities}$$

Therefore we always REJECT (until $P = 1$ exactly)

REASONS :-

- * Risk of classifying $R(\text{classify}) = \lambda_S(1-P) > 0$
for any $P < 1$

- * Risk of rejecting $R(\text{reject}) = 0$

CASE-2 :- $\lambda_R > \lambda_S$ (Rejection expensive)

Ex:- $\lambda_R = 1.5, \lambda_S = 1.0$

$$\tau = 1 - \frac{1.5}{1.0} = -0.5$$

$$P(w_i | x) > -0.5 \text{ (Satisfied)}$$

therefore we NEVER REJECT

REASONS :-

- * Negative threshold means any valid probability exceeds it.
- * Rejection costs more than errors, so always classify.

CASE-3 $\lambda_R < \lambda_S$ (Rejection is cheap)

Ex:- $\lambda_R = 0.5, \lambda_S = 1.0$

$$\tau = 1 - \frac{0.5}{1.0} = 0.5$$

$$P(w_i | x) > 0.5$$

REJECT when max posterior ≤ 0.5 (uncertain)

REASONS :-

- * When uncertain ($P \leq 0.5$), risk of error exceeds cost of rejection

CASE-4: $\lambda_x = \lambda_s$ (Equal costs)

$$\gamma = 1 - \frac{1}{1} = 0$$

$$p(w_i | x) > 0 \quad [\text{Almost always true.}]$$

Therefore we neverly reject (if all posteriors = 0)

Summary

$\lambda_x = 0$ Always reject.

$\lambda_x > \lambda_s \rightarrow$ choose the class with the highest posterior probability as rejection is too costly.

QUESTION-5:

Let Z be drawn from a categorical distribution (takes discrete values) with K possible outcomes/ states and parameter θ , represented by $\text{Cat}(\Theta)$. Describe the value/state using a 1-of- K scheme for $z = [z_1, \dots, z_K]^T$ where $z_k = 1$ if variable is in state k and $z_k = 0$ otherwise. Let the parameter vector for the pdf be $\Theta = [\theta_1, \dots, \theta_K]^T$, where $P(z_k = 1) = \theta_k$, for $k \in \{1, \dots, K\}$.

Given $D\{z_1, \dots, z_N\}$ with iid samples $z_n \sim \text{Cat}(\Theta)$ for $n \in \{1, \dots, N\}$:

Q(5)

A random variable $Z \sim \text{Cat}(\theta)$ with K states

Encoding: 1-of- K representation

$$z = [z_1, \dots, z_K]^T \text{ where } z_k = \begin{cases} 1 & \text{if state } k \\ 0 & \text{otherwise} \end{cases}$$

Given Parameters

$$\theta = [\theta_1, \dots, \theta_K]^T \text{ where } P(z_k = 1) = \theta_k$$

Constraint

$$\sum_{k=1}^K \theta_k = 1$$

Data: $D = \{z_1, \dots, z_N\}$ N i.i.d samples

PART-A

For ML estimator..

For a single sample z_n

$$P(z_n | \theta) = \prod_{k=1}^K \theta_k^{z_{nk}}$$

Exactly one $z_{nk} = 1$ and others are 0

$$\text{If } z_{nk} = 1 \rightarrow \theta_k' = \theta_k$$

$$z_{nk} = 0 \rightarrow \theta_k^0 = 1$$

Product gives θ_k for the active state

For N independent samples

$$P(D | \theta) = \prod_{n=1}^N P(z_n | \theta) = \prod_{n=1}^N \prod_{k=1}^K \theta_k^{z_{nk}}$$

$$P(D | \theta) = \prod_{n=1}^N \prod_{k=1}^K \theta_k^{z_{nk}}$$

Swap the order of products:

$$= \prod_{k=1}^K \prod_{n=1}^N \theta_k^{z_{nk}}$$

$$= \prod_{k=1}^K \theta_k^{\sum_{n=1}^N z_{nk}}$$

Defining $N_k = \sum_{n=1}^N z_{nk}$

$$p(D|\theta) = \prod_{k=1}^K \theta_k^{N_k}$$

Taking the Log-likelihood

$$l(\theta) = \log p(D|\theta) = \log \prod_{k=1}^K \theta_k^{N_k}$$

$$= \sum_{k=1}^K N_k \log \theta_k$$

Let us maximize $l(\theta)$ subject to $\sum_{k=1}^K \theta_k = 1$

Lagrangian:-

$$\mathcal{L}(\theta, \lambda) = \sum_{k=1}^K N_k \log \theta_k + \lambda \left(1 - \sum_{k=1}^K \theta_k \right)$$

Taking partial derivatives w.r. to θ_k .

$$\frac{\partial \mathcal{L}}{\partial \theta_k} = \frac{N_k}{\theta_k} - \lambda = 0$$

Solving for θ_k

$$\frac{N_k}{\theta_k} = \lambda$$

$$\theta_k = \frac{N_k}{\lambda}$$

PART-B

Maximum A Posteriori (MAP) Estimator

Specifying the Prior: Dirichlet Distribution

$$p(\theta | \alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

wherein the normalisation constant is:

$$B(\alpha) = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}$$

Here hyperparameters: $\alpha = [\alpha_1, \dots, \alpha_K]^T$ with $\alpha_k > 0$

Forming the Posterior, using Bayes rule.

$$p(\theta | D) = \frac{p(D | \theta) \cdot p(\theta | \alpha)}{p(D)}$$

To find maximum, we ignore $p(D)$

$$p(\theta | D) \propto p(D | \theta) \cdot p(\theta | \alpha)$$

Now we substitute Likelihood and Prior

$$p(\theta | D) \propto \left[\prod_{k=1}^K \theta_k^{N_k} \right] \cdot \left[\frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \right]$$

Ignore constant $B(\alpha)$:

$$\propto \prod_{k=1}^K \theta_k^{N_k} \cdot \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

$$p(\theta | D) \propto \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1}$$

Log-Posterior

$$\log p(\theta | D) \propto \sum_{k=1}^K (N_k + \alpha_k - 1) \log \theta_k$$

Lagrangian:

$$\mathcal{L}(\theta, \lambda) = \sum_{k=1}^K (N_k + \alpha_k - 1) \log \theta_k + \lambda \left(1 - \sum_{k=1}^K \theta_k \right)$$

Taking Partial Derivatives:

$$\frac{\partial \mathcal{L}}{\partial \theta_k} = \frac{N_k + \alpha_k - 1}{\theta_k} - \lambda = 0$$

Solving, we get

$$\theta_k = \frac{N_k + \alpha_k - 1}{\lambda}$$

Utilizing the constraint to find λ

$$\sum_{k=1}^K \theta_k = 1$$

$$\sum_{k=1}^K \frac{N_k + \alpha_k - 1}{\lambda} = 1$$

$$\frac{1}{\lambda} \sum_{k=1}^K (N_k + \alpha_k - 1) = 1$$

$$\frac{1}{\lambda} \left[\sum_{k=1}^K N_k + \sum_{k=1}^K \alpha_k - K \right] = 1$$

$$\frac{1}{\lambda} \left[N + \sum_{k=1}^K \alpha_k - K \right] = 1$$

$$\lambda = N + \sum_{k=1}^K \alpha_k - K$$

Final MAP Estimator

$$\hat{\theta}_k^{\text{MAP}} = \frac{N_k + \alpha_k - 1}{N + \sum_{j=1}^K \alpha_j - K}$$

Here $\alpha_k - 1$ acts as virtual observation for state k

Push $\alpha_k - 1$ to the count before computing frequency

Special Case Scenario

CASE - 1 :- If $\alpha_k = 1$ for all k

$$\hat{\theta}_k^{\text{MAP}} = \frac{N_k + 1 - 1}{N + K - K} = \frac{N_k}{N} = \hat{\theta}_k^{\text{ML}}$$

MAP reduces to ML \rightarrow Uniform prior

CASE - 2 :- If $\alpha_k = 0.5$ for all k :-

$$\hat{\theta}_k^{\text{MAP}} = \frac{N_k - 0.5}{N - K/2}$$

CASE - 3 :- $N \rightarrow \infty$

$$\hat{\theta}_k^{\text{MAP}} = \frac{N_k + \alpha_k - 1}{N + \sum_j \alpha_j - K} \approx \frac{N_k}{N} \approx \hat{\theta}_k^{\text{ML}}$$

Link for Code Online Repository (Github):

<https://github.com/shre2405/MLPR-Assignment-2>