

Link to code: [Github](#)

Q1

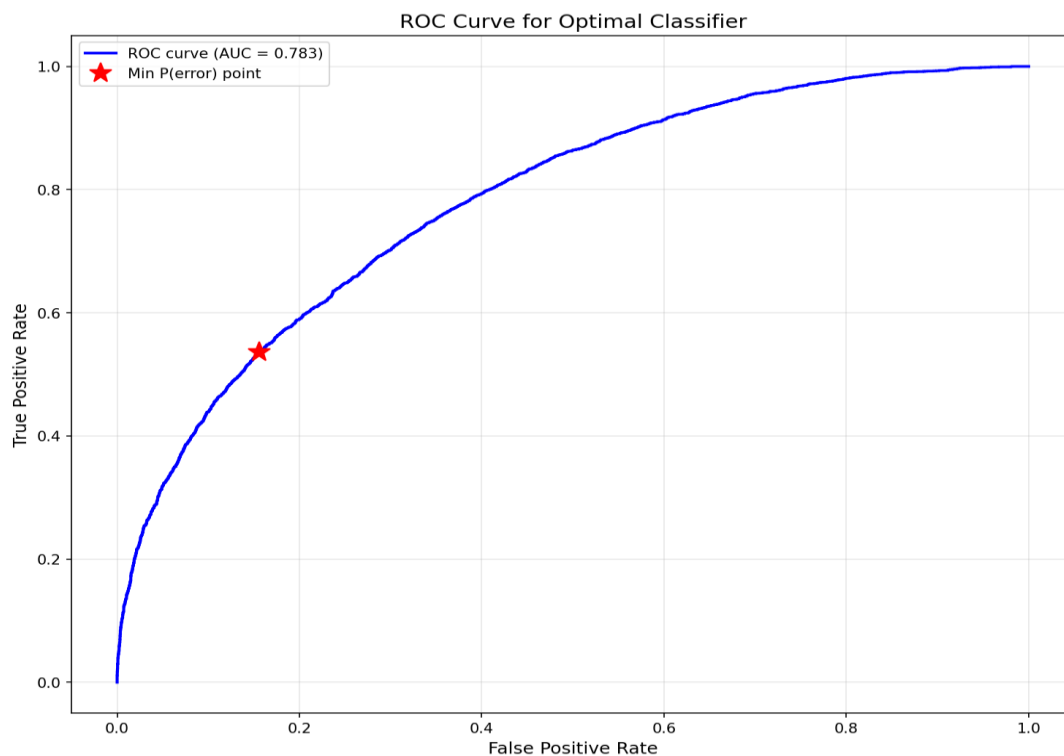
Part 1

$$p(L = 0/x) = \frac{p(x/L = 0)p(L = 0)}{p(x)}$$

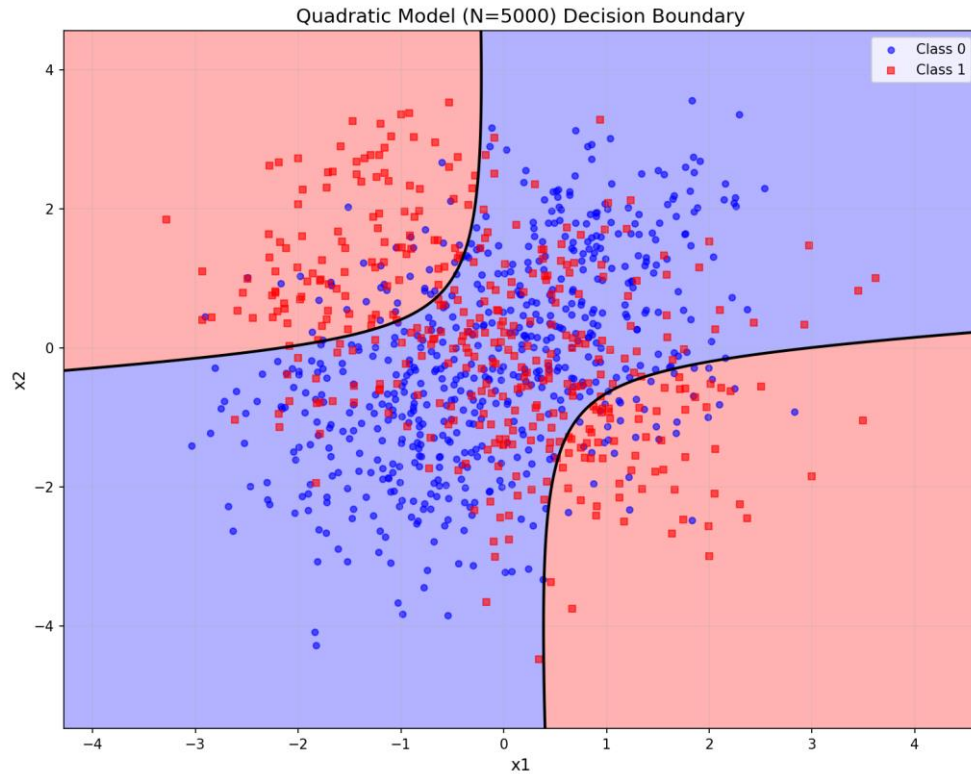
$$p(L = 1/x) = \frac{p(x/L = 1)p(L = 1)}{p(x)}$$

If $p(L = 0/x) > p(L = 1/x) \rightarrow \text{Class 0}$

$p(L = 0/x) < p(L = 1/x) \rightarrow \text{Class 1}$



ROC curve for optimal Bayes classifier on 10,000 validation samples with area under the curve=0.783. The red star marks the min- $P(\text{error})=0.2807$ when $P(L=1|x) = 0.5 = P(L=0|x)$. The linear model performed significantly worse, with an error of 0.5151 at $N=50$ and 0.4047 at $N=5000$, though it is similar to the error of 0.4003 at $N=500$. Quadratic model performed comparably to Bayesian classifier error of 0.2988 at $N=50$ and an error of 0.2855 at $N=5000$. The linear model performs worse due to wrong model selection, whereas the quadratic model captures complexity sufficiently.



From the decision boundary above, there is significant class overlap visible in the central region where the four Gaussian components intersect.

Q2

$$\mathbf{y} = \mathbf{c}(\mathbf{x}, \mathbf{w}) + \mathbf{v} \text{ and } \mathbf{v} \sim N(\mathbf{0}, \sigma^2)$$

$$\begin{aligned} c(\mathbf{x}, \mathbf{w}) &= w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_1 x_2 + w_5 x_2^2 + w_6 x_1^3 + w_7 x_1^2 x_2 + w_8 x_1 x_2^2 + w_9 x_2^3 \\ &= \mathbf{w}^T \mathbf{z}(\mathbf{x}) \end{aligned}$$

$$\text{Where } \mathbf{z}(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2, x_1^3, x_1^2 x_2, x_1 x_2^2, x_2^3]^T$$

For ML estimator

$$\mathbf{D} = [(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)]$$

$$p(\mathbf{D} | \mathbf{w}, \sigma^2) = \prod_{n=1}^N (2\pi\sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2} (\mathbf{y}_n - \mathbf{c}(\mathbf{x}_n, \mathbf{w}))^T \sigma^{-2} (\mathbf{y}_n - \mathbf{c}(\mathbf{x}_n, \mathbf{w}))}$$

$$\mathbf{J}(\mathbf{w}) = \ln(p(\mathbf{D} | \mathbf{w}, \sigma^2)) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (\mathbf{y}_n - \mathbf{c}(\mathbf{x}_n, \mathbf{w}))^T (\mathbf{y}_n - \mathbf{c}(\mathbf{x}_n, \mathbf{w}))$$

$$\mathbf{J}(\mathbf{w}) = (\mathbf{Y} - \mathbf{Z}\mathbf{w})^T (\mathbf{Y} - \mathbf{Z}\mathbf{w})$$

$$\Rightarrow (\mathbf{Y}^T - \mathbf{w}^T \mathbf{Z}^T)(\mathbf{Y} - \mathbf{Z}\mathbf{w}) = \mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{Z}\mathbf{w} - \mathbf{w}^T \mathbf{Z}^T \mathbf{Y} + \mathbf{w}^T \mathbf{Z}^T \mathbf{Z}\mathbf{w}$$

$$\mathbf{J}(\mathbf{w}) = \mathbf{Y}^T \mathbf{Y} - 2\mathbf{w}^T \mathbf{Z}^T \mathbf{Y} + \mathbf{w}^T \mathbf{Z}^T \mathbf{Z}\mathbf{w}$$

$$\frac{\partial J(w)}{\partial w} = -2Z^T Y + 2Z^T Z w = 0$$

$$\hat{w}_{ML} = (Z^T Z)^{-1} Z^T Y$$

For MAP estimation

$$p(w|D) = p(D|w)p(w) \text{ and } w \sim N(0, \gamma I)$$

$$\ln(p(w|D)) = \ln(p(D|w)) + \ln(p(w))$$

$$p(w) = (2\pi)^{-d/2} |\gamma I|^{-1/2} e^{-\frac{1}{2} w^T \gamma^{-1} w}$$

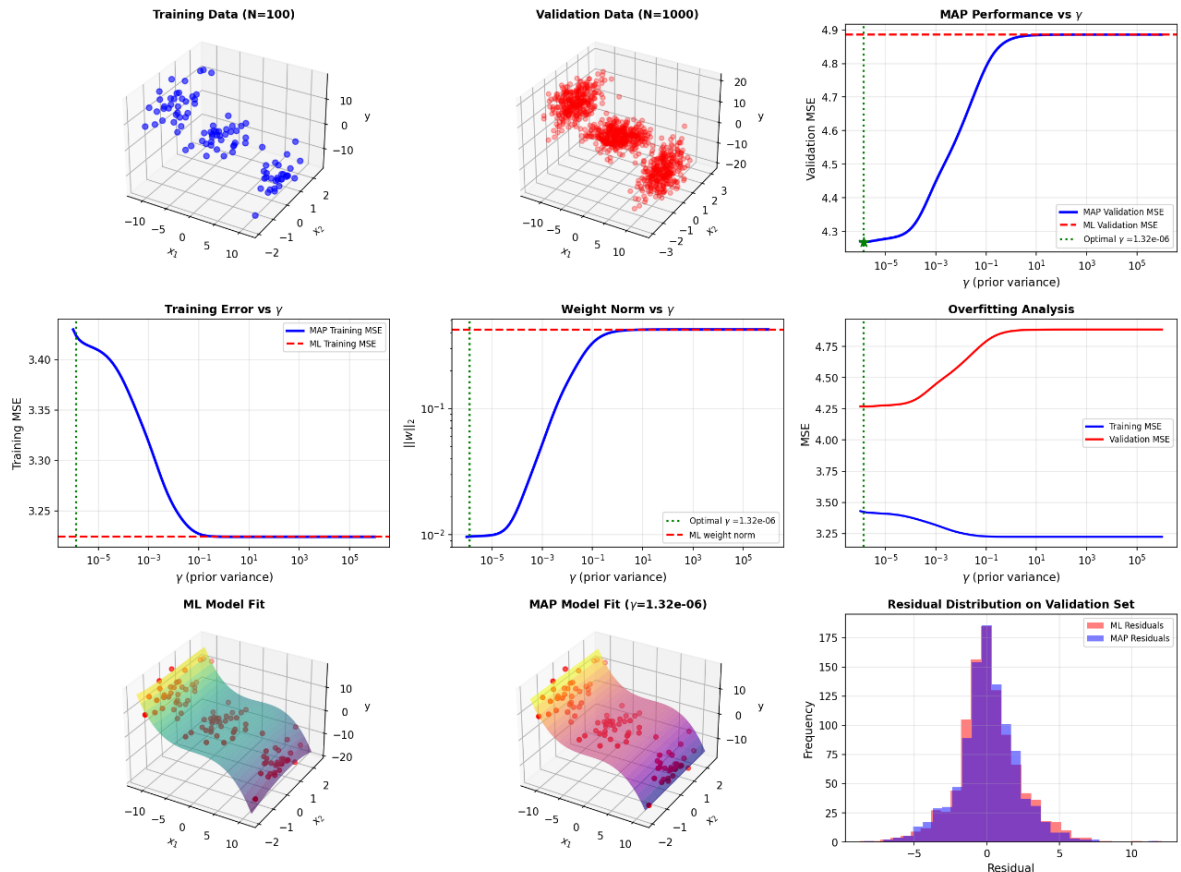
$$\ln(p(w)) = -\frac{1}{2\gamma} w^T w$$

$$\min(J(w)) = -\frac{1}{2\gamma} w^T w - \frac{1}{2\sigma^2} (Y - Zw)^T (Y - Zw)$$

$$\frac{\partial J(w)}{\partial w} = -\frac{2w}{2\gamma} - \frac{2Z^T Z w - 2Z^T Y}{2\sigma^2}$$

$$\hat{w}_{MAP} = \left(\frac{\sigma^2}{\gamma} I + Z^T Z \right)^{-1} Z^T Y$$

We can observe that as $\lim_{\gamma \rightarrow \infty} \hat{w}_{MAP} = \hat{w}_{ML}$



The ML estimator results in overfitting, with a Validation MSE of 4.8862. This is because there are only 100 samples for training and 1000 samples for validation. MAP performs better than the ML estimator due to the regularisation term ($\frac{\sigma^2}{\gamma} \mathbf{I}$). Best γ : 1.3219e-06 and MAP estimator MSE with validation data was 4.2684. Which resulted in an improvement of 12.64% compared to ML estimator

Q3

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

$$\mathbf{p}(\mathbf{r}_i|\mathbf{x}, \mathbf{y}) = (2\pi\sigma_i^2)^{-1/2} \exp \left[-\frac{(\mathbf{r}_i - \mathbf{d}_i)^2}{2\sigma_i^2} \right]$$

$$\text{Where } \mathbf{d}_i = \sqrt{(x - \mathbf{x}_i)^2 + (y - \mathbf{y}_i)^2}$$

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

$$\ln(\mathbf{p}(\mathbf{x}, \mathbf{y})) = -\frac{1}{2} \left\{ \frac{\mathbf{x}^2}{\sigma_x^2} + \frac{\mathbf{y}^2}{\sigma_y^2} \right\}$$

$$\ln(\mathbf{p}(\mathbf{r}|\mathbf{x}, \mathbf{y})) = -\sum \frac{(\mathbf{r}_i - \mathbf{d}_i)^2}{2\sigma_i^2}$$

$$\ln(\mathbf{p}(\mathbf{x}, \mathbf{y}|\mathbf{r})) = -\sum \frac{(\mathbf{r}_i - \mathbf{d}_i)^2}{2\sigma_i^2} - \frac{1}{2} \left[\frac{\mathbf{x}^2}{\sigma_x^2} + \frac{\mathbf{y}^2}{\sigma_y^2} \right]$$

For MAP estimation $[\mathbf{x}_{\text{MAP}}, \mathbf{y}_{\text{MAP}}]^T = \min(\ln(\mathbf{p}(\mathbf{x}, \mathbf{y}|\mathbf{r})))$

$$= \min \left(-\sum \frac{(\mathbf{r}_i - \mathbf{d}_i)^2}{2\sigma_i^2} - \frac{1}{2} \left[\frac{\mathbf{x}^2}{\sigma_x^2} + \frac{\mathbf{y}^2}{\sigma_y^2} \right] \right)$$

The true position of the vehicle is set at $[0.3, 0.4]^T$, noise variance of $\sigma_i = 0.3$ and $\sigma_x = \sigma_y = 0.25$. The code generates evenly spaced K landmarks on a unit circle. Noise data was generated. Then the log likelihood is set up according to the equation above. To get $[\mathbf{x}_{\text{MAP}}, \mathbf{y}_{\text{MAP}}]^T$ we have to minimize the above equation.

K = 1 Landmarks (on unit circle): Landmark 1: [1.000, 0.000]

r1 = 0.9552 (true distance = 0.8062, noise_error = +0.1490)

MAP estimate: [0.0183, -0.0000], Localization error: 0.4892

K = 2 Landmarks (on unit circle): Landmark 1: [1.000, 0.000] Landmark 2: [-1.000, 0.000]

r1 = 0.7647 (true distance = 0.8062, noise_error = -0.0415)

r2 = 1.5545 (true distance = 1.3601, noise_error = +0.1943)

MAP estimate: [0.2296, -0.0000], Localization error: 0.4062

K = 3 Landmarks (on unit circle): Landmark 1: [1.000, 0.000] Landmark 2: [-0.500, 0.866] Landmark 3: [-0.500,-0.866]

r1 = 1.2631 (true distance = 0.8062, noise_error = +0.4569)

r2 = 0.8556 (true distance = 0.9258, noise_error= -0.0702)

r3 = 1.4274 (true distance = 1.4976, noise_error = -0.0702)

MAP estimate: [-0.0457, 0.1840], Localization error: 0.4077

K = 4 Landmarks (on unit circle): Landmark 1: [1.000, 0.000] Landmark 2: [0.000, 1.000] Landmark 3: [-1.000, 0.000] Landmark 4: [-0.000,-1.000]

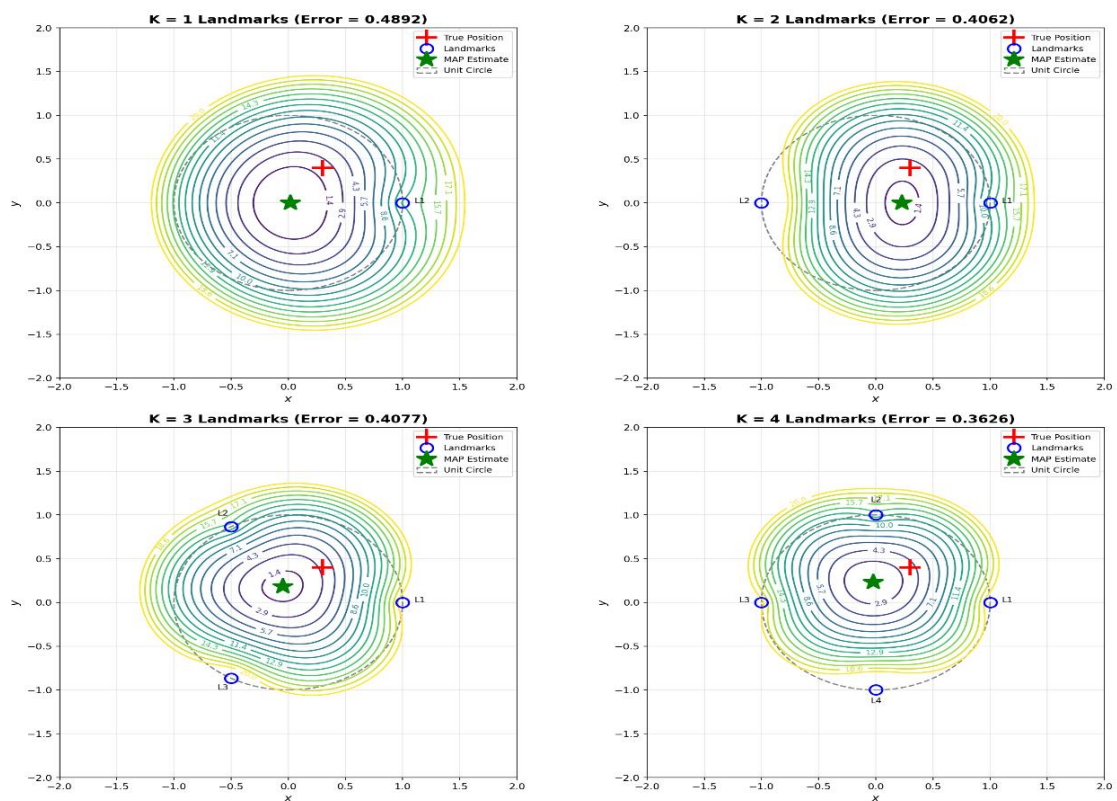
r1 = 1.2800 (true distance = 0.8062, noise_error= +0.4738)

r2 = 0.9011 (true distance = 0.6708, noise_error = +0.2302)

r3 = 1.2193 (true distance = 1.3601, noise_error = -0.1408)

r4 = 1.5946 (true distance = 1.4318, noise_error = +0.1628)

MAP estimate: [-0.0207, 0.2308], Localization error: 0.3626



We observe that as the number of K increases, localization error decreases. As K increases, the contours gradually become circle-shaped, centred around the estimated position. The numbers represent the value of the objective function for each (x,y) value. For K=1, the innermost visible contours (J = 2.9) cover a large area, indicating the minimum is spread over a broad, shallow region. For K=4, the innermost contours (J \approx 2.9) cover a much smaller area, which means higher certainty.

Q4

$$\lambda(\alpha_i, w_j) = \begin{cases} 0, & i = j = 1, 2, \dots, c \\ \lambda_r, & c + 1 \\ \lambda_s, & \text{otherwise} \end{cases}$$

$$R(\alpha_i, x) = \sum_{j=1}^c \lambda_s(\alpha_i / w_j) p(w_j / x)$$

$$= \lambda_s \sum_{i \neq j}^c p(w_j / x) = \lambda_s(1 - p(w_i / x))$$

$$R(\text{rejection} / x) = \lambda_r$$

If we decide that w_i class is better than all the other classes, the risk for class w_i should be the least among all classes and risk of w_i should be less than rejection

Condition 1. $R(\alpha_i / x) \leq R(\alpha_j / x)$

$$\lambda_s(1 - p(w_i / x)) \leq \lambda_s(1 - p(w_j / x))$$

$$p(w_i / x) \geq p(w_j / x)$$

and **Condition 2.** $R(\alpha_i / x) < R(\text{rejection} / x)$

$$\lambda_s(1 - p(w_i / x)) < \lambda_r$$

$$p(w_i / x) > 1 - \frac{\lambda_r}{\lambda_s}$$

Otherwise reject

Special Case 1.

$$\lambda_r = 0$$

This makes rejection free of loss which in turn makes **Condition2** $p(w_i / x) > 1$ which is not possible as all probabilities are ≤ 1 . This makes the classifier always rejecting

Special Case 2.

$$\lambda_r > \lambda_s$$

$$1 - \frac{\lambda_r}{\lambda_s} < 0$$

Condition 2 now becomes always satisfied as all probabilities are between 0 and 1. This makes classifier to never reject instead always classifies data to most likely class

Q5

ML estimation for Θ

$$p(\mathbf{D}|\Theta) = \prod_{n=1}^N p(z_n|\theta)$$

$$p(z|\Theta) = \prod_{k=1}^K \theta_k^{z_k}$$

$$p(\mathbf{D}|\Theta) = \prod_{k=1}^K \prod_{n=1}^N \theta_k^{z_{nk}} = \prod_{k=1}^K \theta_k^{N_k} \text{ here } N_k = \sum_{n=1}^N z_{nk}$$

$$\ln(p(\mathbf{D}|\Theta)) = \sum_{k=1}^K N_k \ln(\theta_k) \text{ and constraint } \sum_{k=1}^K \theta_k = 1$$

Using Lagrange multiplier λ

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

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

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

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

$$\hat{\theta}_{ML} = \frac{N_k}{\sum_{k=1}^K N_k}$$

If $p(\Theta)$ has a Dirichlet distribution

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

$$p(\Theta|\mathbf{D}, \alpha) = p(\mathbf{D}|\Theta, \alpha)p(\Theta|\alpha)$$

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

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

$$\ln(\mathbf{p}(\boldsymbol{\Theta}|\mathbf{D}, \alpha)) = \sum_{k=1}^K (N_k + \alpha_k - 1) \ln(\boldsymbol{\theta}_k) \text{ and constraint } \sum_{k=1}^K \boldsymbol{\theta}_k = \mathbf{1}$$

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

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}_k} = \frac{(N_k + \alpha_k - 1)}{\boldsymbol{\theta}_k} - \lambda = 0$$

$$\boldsymbol{\theta}_k = \frac{(N_k + \alpha_k - 1)}{\lambda}$$

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

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

$$\hat{\boldsymbol{\theta}}_{MAP} = \frac{(N_k + \alpha_k - 1)}{\sum_{k=1}^K (N_k + \alpha_k - 1)}$$

References

<https://scikit-learn.org/stable/modules/preprocessing.html#polynomial-features>