

Lecture 11, Feb 12, 2024

Gaussian Discriminant Analysis

- Consider a classification problem where we have classes $c \in C$, each having a prior π_c , jointly Gaussian distributed with a mean of $\boldsymbol{\mu}_c$ and a covariance of $\boldsymbol{\Sigma}_c$
- Gaussian discriminant analysis is a special case of hypothesis testing for this type of classification problem; given an observation of the vector \mathbf{X} , we would like to know which class it came from (i.e. which hypothesis is true)
- Consider the case where all the classes have the same covariance $\boldsymbol{\Sigma}$, so they differ only by their mean
 - The posterior is $P[y = c|\mathbf{x}] = \frac{f(\mathbf{x}|c)\pi_c}{\sum_{c'} f(\mathbf{x}|c')\pi_{c'}}$
 - The numerator becomes $\frac{e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_c)^T\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}_c)}}{(2\pi)^{\frac{D}{2}}\sqrt{\det \boldsymbol{\Sigma}}}\pi_c = \frac{e^{-\frac{1}{2}\mathbf{x}^T\boldsymbol{\Sigma}^{-1}\mathbf{x}}}{(2\pi)^{\frac{D}{2}}\sqrt{\det \boldsymbol{\Sigma}}}\exp\left[-\frac{1}{2}\boldsymbol{\mu}_c^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_c + \boldsymbol{\mu}_c^T\boldsymbol{\Sigma}^{-1}\mathbf{x} + \log \pi_c\right]$
 - Let $\boldsymbol{\beta}_c^T = \boldsymbol{\mu}_c^T\boldsymbol{\Sigma}^{-1}$ and $\gamma_c = \log \pi_c - \frac{1}{2}\boldsymbol{\mu}_c^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_c$
 - The exponential can then be written as $p(y = c|\mathbf{x}) = \frac{\exp(\boldsymbol{\beta}_c^T\mathbf{x} + \gamma_c)}{\sum_{c'} \exp(\boldsymbol{\beta}_{c'}^T\mathbf{x} + \gamma_{c'})}$
 - * This is a softmax function
 - * The exponential in the softmax makes it so that the largest term dominates while all other terms are usually much smaller
 - * Each class has an associated $\boldsymbol{\beta}_c$ and γ_c , which contains all the info of the class
 - With this, we have $p(y = c|x) \approx \begin{cases} 1 & \boldsymbol{\beta}_c^T\mathbf{x} + \gamma_c \gg \boldsymbol{\beta}_{c'}^T\mathbf{x} + \gamma_{c'} \\ 0 & \text{otherwise} \end{cases}$
 - The decision rule is $\hat{y}(\mathbf{x}) = \underset{c}{\operatorname{argmax}} \boldsymbol{\beta}_c^T\mathbf{x} + \gamma_c$
 - This is referred to as *linear Gaussian discriminant analysis*, since the decision boundary is a linear function of \mathbf{x}
 - * The boundary occurs where $\boldsymbol{\beta}_1^T\mathbf{x} + \gamma_1 = \boldsymbol{\beta}_0^T\mathbf{x} + \gamma_0$ which forms a hyperplane
- More generally, the covariances of the classes are different
 - $\log p(y = c|\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_c)^T\boldsymbol{\Sigma}_c^{-1}(\mathbf{x} - \boldsymbol{\mu}_c) - \frac{1}{2}\log \det \boldsymbol{\Sigma}_c + \log \pi_c - \frac{D}{2}\log 2\pi$
 - Consider the boundary between two regions:
 - * $-\frac{1}{2}\mathbf{x}^T\boldsymbol{\Sigma}_1^{-1}\mathbf{x} + \boldsymbol{\mu}_1^T\boldsymbol{\Sigma}_1^{-1}\mathbf{x} - \frac{1}{2}\log \det \boldsymbol{\Sigma}_1 + \log \pi_1 = -\frac{1}{2}\mathbf{x}^T\boldsymbol{\Sigma}_0^{-1}\mathbf{x} + \boldsymbol{\mu}_0^T\boldsymbol{\Sigma}_0^{-1}\mathbf{x} - \frac{1}{2}\log \det \boldsymbol{\Sigma}_0 + \log \pi_0$
 - * $\mathbf{x}^T(\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_0^{-1})\mathbf{x} + 2(\boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0 - \boldsymbol{\Sigma}_1^{-1}\boldsymbol{\mu}_1)^T\mathbf{x} + (\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1^T\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\mu}_1) + \log \frac{\det \boldsymbol{\Sigma}_0}{\det \boldsymbol{\Sigma}_1} + 2\log \frac{\pi_1}{\pi_0} = 0$
 - * This is a quadratic form that can define a parabola, hyperbola, or even circles and ellipses
- To obtain the parameters of the Gaussian distribution of each class, we can use ML estimation
 - $\hat{\pi}_c = \frac{n_c}{n}$ is given by the relative frequency of class c
 - $\hat{\boldsymbol{\mu}}_c = \frac{1}{n_c} \sum_i \mathbf{x}_i^c$ is given by the sample mean
 - $\hat{\boldsymbol{\Sigma}}_c = \frac{1}{n_c} (\mathbf{x}_i^c - \hat{\boldsymbol{\mu}}_c)^T (\mathbf{x}_i^c - \hat{\boldsymbol{\mu}}_c)$
 - Note that the variance estimate here is biased; use the version with $n_c - 1$ in the denominator for unbiased

Gaussian Parameter Estimation

- Let \mathbf{X}, \mathbf{Y} be jointly Gaussian and let $\mathbf{w} = \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}$; we want to find the MAP estimator for \mathbf{X} given \mathbf{Y}
 - i.e. we want to find the distribution of \mathbf{X} conditioned on \mathbf{Y}
 - We will make use of the covariances between elements of \mathbf{X} and \mathbf{Y}
 - Strategy: expand out the exponent of the joint PDF and rearrange it into a form with \mathbf{X} as the variable

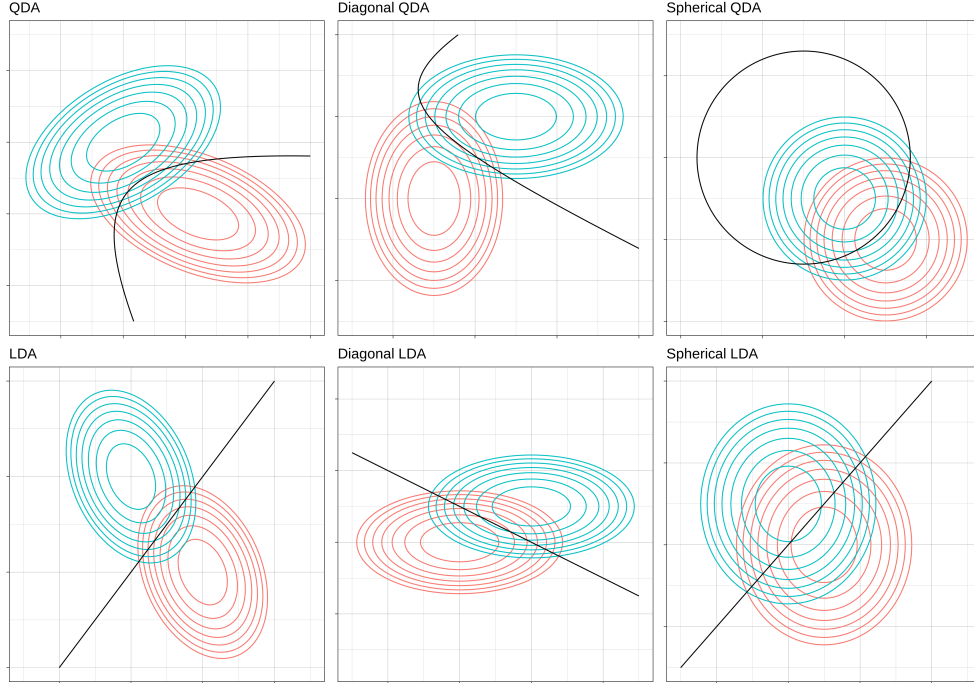


Figure 1: Illustration of the different kinds of Gaussian discriminant analysis.

- The mean of \mathbf{w} is $\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_Y \end{bmatrix}$
- The covariance is $\boldsymbol{\Sigma}_w = E[(\mathbf{w} - \boldsymbol{\mu}_w)(\mathbf{w} - \boldsymbol{\mu}_w)^T] = E \left[\begin{bmatrix} \mathbf{X} - \boldsymbol{\mu}_X \\ \mathbf{Y} - \boldsymbol{\mu}_Y \end{bmatrix} \begin{bmatrix} (\mathbf{X} - \boldsymbol{\mu}_X)^T \\ (\mathbf{Y} - \boldsymbol{\mu}_Y)^T \end{bmatrix} \right] = \begin{bmatrix} \boldsymbol{\Sigma}_{XX} & \boldsymbol{\Sigma}_{XY} \\ \boldsymbol{\Sigma}_{YX} & \boldsymbol{\Sigma}_{YY} \end{bmatrix}$
 - The overall dimension is $N \times N$; the $\boldsymbol{\Sigma}_{XY}, \boldsymbol{\Sigma}_{YX}$ matrices are in general rectangular
 - Let the *precision matrix* $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}_w^{-1} = \begin{bmatrix} \boldsymbol{\Lambda}_{XX} & \boldsymbol{\Lambda}_{XY} \\ \boldsymbol{\Lambda}_{YX} & \boldsymbol{\Lambda}_{YY} \end{bmatrix}$
 - * This is the opposite of variance; the larger the precision, the more tightly concentrated the distribution
 - * Note $\boldsymbol{\Sigma}_{XX} \neq \boldsymbol{\Lambda}_{XX}^{-1}$, and $\boldsymbol{\Lambda}_{XY} = \boldsymbol{\Lambda}_{YX}^T$
- Now consider the exponent of the joint PDF
 - $-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_w)^T \boldsymbol{\Sigma}_w^{-1} (\mathbf{w} - \boldsymbol{\mu}_w)$

$$= -\frac{1}{2} \begin{bmatrix} \mathbf{X} - \boldsymbol{\mu}_X \\ \mathbf{Y} - \boldsymbol{\mu}_Y \end{bmatrix}^T \begin{bmatrix} \boldsymbol{\Lambda}_{XX} & \boldsymbol{\Lambda}_{XY} \\ \boldsymbol{\Lambda}_{YX} & \boldsymbol{\Lambda}_{YY} \end{bmatrix} \begin{bmatrix} \mathbf{X} - \boldsymbol{\mu}_X \\ \mathbf{Y} - \boldsymbol{\mu}_Y \end{bmatrix}$$

$$= -\frac{1}{2}(\mathbf{X} - \boldsymbol{\mu}_X)^T \boldsymbol{\Lambda}_{XX}(\mathbf{X} - \boldsymbol{\mu}_X) - \frac{1}{2}(\mathbf{X} - \boldsymbol{\mu}_X)^T \boldsymbol{\Lambda}_{XY}(\mathbf{Y} - \boldsymbol{\mu}_Y) - \frac{1}{2}(\mathbf{Y} - \boldsymbol{\mu}_Y)^T \boldsymbol{\Lambda}_{YX}(\mathbf{X} - \boldsymbol{\mu}_X) - \frac{1}{2}(\mathbf{Y} - \boldsymbol{\mu}_Y)^T \boldsymbol{\Lambda}_{YY}(\mathbf{Y} - \boldsymbol{\mu}_Y)$$

$$= -\frac{1}{2}(\mathbf{X}^T \boldsymbol{\Lambda}_{XX} \mathbf{X} - \mathbf{X}^T \boldsymbol{\Lambda}_{XY}(\mathbf{Y} - \boldsymbol{\mu}_Y) + (\mathbf{Y} - \boldsymbol{\mu}_Y)^T \boldsymbol{\Lambda}_{YX} \mathbf{X} - \boldsymbol{\mu}_X^T \boldsymbol{\Lambda}_{XX} \mathbf{X} - \mathbf{X}^T \boldsymbol{\Lambda}_{XX} \boldsymbol{\mu}_X + \dots)$$

$$= -\frac{1}{2}(\mathbf{X}^T \boldsymbol{\Lambda}_{XX} \mathbf{X} - 2\mathbf{X}^T(\boldsymbol{\Lambda}_{XX} \boldsymbol{\mu}_X - \boldsymbol{\Lambda}_{XY}(\mathbf{Y} - \boldsymbol{\mu}_Y)) + \dots)$$
 - This gives us $f(\mathbf{x}|\mathbf{y})$ and implies that it is jointly Gaussian
 - Let $f(\mathbf{x}|\mathbf{y}) = c \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{X|Y})^T \boldsymbol{\Sigma}_{X|Y}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{X|Y}) \right] = c \exp \left[-\frac{1}{2}(\mathbf{x} \boldsymbol{\sigma}_{XY}^{-1} \mathbf{x} - 2\mathbf{x}^T \boldsymbol{\Sigma}_{XY}^{-1} \boldsymbol{\mu}_{X|Y}) \right]$
 - By matching terms we see $\boldsymbol{\Sigma}_{XY} = \boldsymbol{\Lambda}_{XX}^{-1}$ and $\boldsymbol{\mu}_{X|Y} = \boldsymbol{\mu}_X - \boldsymbol{\Lambda}_{XX}^{-1} \boldsymbol{\Lambda}_{XY}(\mathbf{y} - \boldsymbol{\mu}_Y)$
 - Therefore $f_{X|Y}(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}_{X|Y}, \boldsymbol{\Sigma}_{X|Y})$
- Given this PDF, we can see that the MAP estimate is simply $\boldsymbol{\mu}_{X|Y}$
 - We can show that this is the same as the LMS estimate

- However, we only have this in terms of the precision matrix; can we find it in terms of Σ ?
- Using the Schur complement on Σ^{-1} we can find a general expression for each of the Λ
 - $\Lambda_{XX} = (\Sigma_{XX} - \Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{YX})^{-1}$
 - $\Lambda_{XY} = (\Sigma_{XX} - \Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{YX})^{-1}\Sigma_{XY}\Sigma_{YY}^{-1}$
 - Therefore $\mu_{X|Y} = \mu_X + \Sigma_{XY}\Sigma_{YY}^{-1}(y - \mu_Y)$, $\Sigma_{X|Y} = \Sigma_{XX} - \Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{YX}$
 - We get both the mean of the estimate and its spread