# Linear and Quadratic Discriminant Analysis

Statistical Machine Learning (ENGG*6600*08)
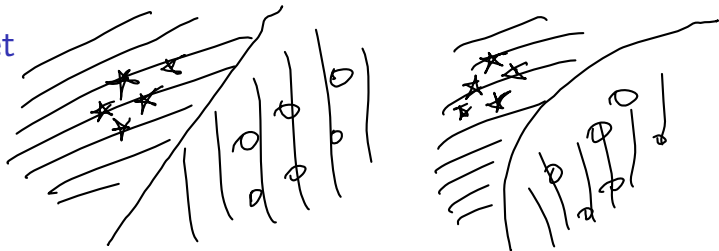
School of Engineering,
University of Guelph, ON, Canada

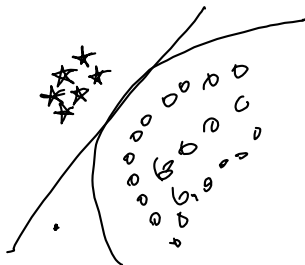Course Instructor: Benyamin Ghojogh
Fall 2023

**Optimization for the Boundary of Classes**

# Dataset



- Assume we have a dataset of *instances* $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$ with sample size $n$ and dimensionality $\boldsymbol{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$. The $y_i$'s are the class labels.
- We would like to *classify* the space of data using these instances.
- Linear Discriminant Analysis (LDA) and Quadratic discriminant Analysis (QDA) [1] are two well-known *supervised* classification methods in statistical and probabilistic learning.

# Optimization for the Boundary of Classes

- First suppose the data is one dimensional, $x \in \mathbb{R}$. Assume we have two classes with the Cumulative Distribution Functions (CDF) $F_1(x)$ and $F_2(x)$, respectively. Let the Probability Density Functions (PDF) of these CDFs be:

$$\begin{cases} f_1(x) = \dfrac{\partial F_1(x)}{\partial x}, & (1) \\[2mm] f_2(x) = \dfrac{\partial F_2(x)}{\partial x}, & (2) \end{cases}$$
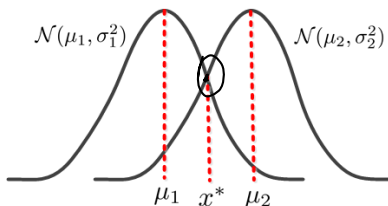
  respectively.

- We assume that the two classes have normal (Gaussian) distribution which is the most common and default distribution in the real-world applications. The mean of one of the two classes is greater than the other one; we assume $\mu_1 < \mu_2$. An instance $x \in \mathbb{R}$ belongs to one of these two classes:

$$x \sim \begin{cases} \mathcal{N}(\mu_1, \sigma_1^2), & \text{if } x \in \mathcal{C}_1, \\ \mathcal{N}(\mu_2, \sigma_2^2), & \text{if } x \in \mathcal{C}_2, \end{cases} \tag{3}$$

  where $\mathcal{C}_1$ and $\mathcal{C}_2$ denote the first and second class, respectively.

# Optimization for the Boundary of Classes

- For an instance $x$, we may have an error in estimation of the class it belongs to. At a point, which we denote by $x^*$, the probability of the two classes are equal; therefore, the point $x^*$ is on the boundary of the two classes.

- As we have $\mu_1 < \mu_2$, we can say $\mu_1 < x^* < \mu_2$ as shown in below figure. Therefore, if $x < x^*$ or $x > x^*$ the instance $x$ belongs to the first and second class, respectively. Hence, estimating $x < x^*$ or $x > x^*$ for belonging to the second and first class, respectively, is an error in estimation of the class.

# Optimization for the Boundary of Classes

- This probability of the error can be stated as:

$$\mathbb{P}(\text{error}) = \mathbb{P}(x > x^*, x \in \mathcal{C}_1) + \mathbb{P}(x < x^*, x \in \mathcal{C}_2). \tag{4}$$
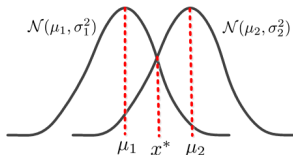
- As we have $\mathbb{P}(A, B) = \mathbb{P}(A|B)\, \mathbb{P}(B)$, we can say:

likelihood      prior

$$\mathbb{P}(\text{error}) = \mathbb{P}(x > x^* \mid x \in \mathcal{C}_1)\, \mathbb{P}(x \in \mathcal{C}_1) + \mathbb{P}(x < x^* \mid x \in \mathcal{C}_2)\, \mathbb{P}(x \in \mathcal{C}_2), \tag{5}$$

which we want to minimize:

$$\underset{x^*}{\text{minimize}} \;\; \mathbb{P}(\text{error}), \tag{6}$$

by finding the best boundary of classes, i.e., $x^*$.

# Optimization for the Boundary of Classes

- We found:

$$\mathbb{P}(\text{error}) = \mathbb{P}(x > x^* \mid x \in \mathcal{C}_1)\mathbb{P}(x \in \mathcal{C}_1) + \mathbb{P}(x < x^* \mid x \in \mathcal{C}_2)\mathbb{P}(x \in \mathcal{C}_2).$$

- According to the definition of CDF, we have:

$$\mathbb{P}(x < c, x \in \mathcal{C}_1) = F_1(c) \implies \mathbb{P}(x > x^*, x \in \mathcal{C}_1) = 1 - F_1(x^*), \tag{7}$$

$$\mathbb{P}(x < x^*, x \in \mathcal{C}_2) = F_2(x^*). \tag{8}$$

$$\mathbb{P}(x < x^*, x \in \mathcal{C}_1) = F_1(x^*)$$

- According to the definition of PDF, we have:

$$\mathbb{P}(x \in \mathcal{C}_1) = f_1(x) = \pi_1, \tag{9}$$

$$\mathbb{P}(x \in \mathcal{C}_2) = f_2(x) = \pi_2, \tag{10}$$

  where we denote the priors $f_1(x)$ and $f_2(x)$ by $\pi_1$ and $\pi_2$, respectively.

- Hence, Eqs. (5) and (6) become:

$$\underset{x^*}{\text{minimize}} \quad (1 - F_1(x^*))\,\pi_1 + F_2(x^*)\,\pi_2. \tag{11}$$

- We take derivative for the sake of minimization:

$$\frac{\partial\, \mathbb{P}(\text{error})}{\partial x^*} = -f_1(x^*)\,\pi_1 + f_2(x^*)\,\pi_2 \overset{\text{set}}{=} 0 \implies f_1(x^*)\,\pi_1 = f_2(x^*)\,\pi_2. \tag{12}$$

# Optimization for the Boundary of Classes

- Another way to obtain this expression is equating the posterior probabilities to have the equation of the boundary of classes:

$$\mathbb{P}(x \in \mathcal{C}_1 \mid X = x) \overset{\text{set}}{=} \mathbb{P}(x \in \mathcal{C}_2 \mid X = x). \tag{13}$$

- According to Bayes rule, the *posterior* is:

$$\mathbb{P}(x \in \mathcal{C}_1 \mid X = x) = \frac{\mathbb{P}(X = x \mid x \in \mathcal{C}_1)\,\mathbb{P}(x \in \mathcal{C}_1)}{\mathbb{P}(X = x)}$$
$$= \frac{f_1(x)\,\pi_1}{\sum_{k=1}^{|\mathcal{C}|} \mathbb{P}(X = x \mid x \in \mathcal{C}_k)\,\pi_k}, \tag{14}$$

where $|\mathcal{C}|$ is the number of classes which is two here. The $f_1(x)$ and $\pi_1$ are the *likelihood (class conditional)* and *prior* probabilities, respectively, and the denominator is the marginal probability.

- Therefore, Eq. (13) becomes:

$$\frac{f_1(x)\,\pi_1}{\sum_{i=1}^{|\mathcal{C}|} \mathbb{P}(X = x \mid x \in \mathcal{C}_i)\,\pi_i} \overset{\text{set}}{=} \frac{f_2(x)\,\pi_2}{\sum_{i=1}^{|\mathcal{C}|} \mathbb{P}(X = x \mid x \in \mathcal{C}_i)\,\pi_i} \implies f_1(x)\,\pi_1 = f_2(x)\,\pi_2. \tag{15}$$

Bayes' rule:

$$P(\underbrace{x \in c}_{\text{posterior}} | \textcircled{X}) = \frac{\overbrace{P(X | \textcircled{c})}^{\text{likelihood}} \overbrace{P(c)}^{\text{prior \& class}}}{\underbrace{P(X)}_{\substack{\text{prior \& data} \\ (\text{marginalized prob})}}}$$

$$P(X) = \sum_{c} P(X|c) \, P(c)$$

# Optimization for the Boundary of Classes

- Now let us think of data as *multivariate* data with dimensionality $d$. The PDF for multivariate Gaussian distribution, $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is:

$$f(\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp\left(- \frac{(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})}{2}\right), \tag{16}$$

where $\boldsymbol{x} \in \mathbb{R}^d$, $\boldsymbol{\mu} \in \mathbb{R}^d$ is the mean, $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$ is the covariance matrix, and $|.|$ is the determinant of matrix. The $\pi \approx 3.14$ in this equation should not be confused with the $\pi_k$ (prior) in Eq. (12) or (15).

- Therefore, the Eq. (12) or (15) becomes:

$$\begin{cases} \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_1|}} \exp\left(- \frac{(\boldsymbol{x} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_1^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_1)}{2}\right) \pi_1 \\ = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_2|}} \exp\left(- \frac{(\boldsymbol{x} - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_2)}{2}\right) \pi_2, \end{cases} \tag{17}$$

where the distributions of the first and second class are $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$, respectively.

**Linear Discriminant Analysis**

# Linear Discriminant Analysis for Binary Classification

- In Linear Discriminant Analysis (LDA), we assume that the two classes have equal covariance matrices:

$$\mathbf{\Sigma}_1 = \mathbf{\Sigma}_2 = \mathbf{\Sigma}. \tag{18}$$

Therefore, the Eq. (17) becomes:

$$\frac{1}{\sqrt{(2\pi)^d |\mathbf{\Sigma}|}} \exp\left(-\frac{(\boldsymbol{x}-\boldsymbol{\mu}_1)^\top \mathbf{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu}_1)}{2}\right) \pi_1$$

$$= \frac{1}{\sqrt{(2\pi)^d |\mathbf{\Sigma}|}} \exp\left(-\frac{(\boldsymbol{x}-\boldsymbol{\mu}_2)^\top \mathbf{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu}_2)}{2}\right) \pi_2,$$

$$\implies \exp\left(-\frac{(\boldsymbol{x}-\boldsymbol{\mu}_1)^\top \mathbf{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu}_1)}{2}\right) \pi_1 = \exp\left(-\frac{(\boldsymbol{x}-\boldsymbol{\mu}_2)^\top \mathbf{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu}_2)}{2}\right) \pi_2,$$

$$\stackrel{(a)}{\implies} -\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu}_1)^\top \mathbf{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu}_1) + \ln(\pi_1) = -\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu}_2)^\top \mathbf{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu}_2) + \ln(\pi_2),$$

where (a) takes natural logarithm from the sides of equation.

# Linear Discriminant Analysis for Binary Classification

- We can simplify this term as:

$$\underbrace{(x - \mu_1)^\top \Sigma^{-1}(x - \mu_1)} = (x^\top - \mu_1^\top)\Sigma^{-1}(x - \mu_1)$$
$$= x^\top \Sigma^{-1} x - x^\top \Sigma^{-1}\mu_1 - \mu_1^\top \Sigma^{-1} x + \mu_1^\top \Sigma^{-1}\mu_1$$
$$\overset{(a)}{=} x^\top \Sigma^{-1} x + \mu_1^\top \Sigma^{-1}\mu_1 - 2\mu_1^\top \Sigma^{-1} x, \tag{19}$$

where (a) is because $x^\top \Sigma^{-1}\mu_1 = \mu_1^\top \Sigma^{-1} x$ as it is a scalar and $\Sigma^{-1}$ is symmetric so $\Sigma^{-\top} = \Sigma^{-1}$. Thus, we have:

$$-\frac{1}{2}x^\top \Sigma^{-1} x - \frac{1}{2}\mu_1^\top \Sigma^{-1}\mu_1 + \mu_1^\top \Sigma^{-1} x + \ln(\pi_1)$$
$$= -\frac{1}{2}x^\top \Sigma^{-1} x - \frac{1}{2}\mu_2^\top \Sigma^{-1}\mu_2 + \mu_2^\top \Sigma^{-1} x + \ln(\pi_2).$$

- Therefore, if we multiply the sides of equation by 2, we have:

$$\left(2\Sigma^{-1}(\mu_2 - \mu_1)\right)^\top x + \left(\mu_1^\top \Sigma^{-1}\mu_1 - \mu_2^\top \Sigma^{-1}\mu_2\right) + 2\ln\left(\frac{\pi_2}{\pi_1}\right) = 0, \tag{20}$$

which is the equation of a line in the form of $a^\top x + b = 0$.

- Therefore, if we consider Gaussian distributions for the two classes where the covariance matrices are assumed to be equal, the decision boundary of classification is a line. Because of linearity of the decision boundary which discriminates the two classes, this method is named *linear discriminant* analysis.

$$\left(x^T \Sigma^{-1} \mu_1\right) = \left(x^T \Sigma^{-1} \mu_1\right)^T = \mu_1^T \underbrace{\left(\Sigma^{-T}\right.}_{\Sigma^{-1}} \underbrace{\left(x^T\right)^T}_{x}$$

$$= \mu_1^T \Sigma^{-1} x$$

# Linear Discriminant Analysis for Binary Classification

- For obtaining Eq. (20), we brought the expressions to the right side which was corresponding to the second class; therefore, if we use $\delta(\boldsymbol{x}) : \mathbb{R}^d \to \mathbb{R}$ as the left-hand-side expression (function) in Eq. (20):

$$\delta(\boldsymbol{x}) := 2\left(\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)\right)^{\top}\boldsymbol{x} + \boldsymbol{\mu}_1^{\top}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^{\top}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_2 + 2\ln(\frac{\pi_2}{\pi_1}), \tag{21}$$
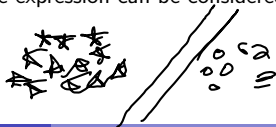
the class of an instance $\boldsymbol{x}$ is estimated as:

$$\widehat{\mathcal{C}}(x) = \begin{cases} 1, & \text{if } \delta(\boldsymbol{x}) < 0, \\ 2, & \text{if } \delta(\boldsymbol{x}) > 0. \end{cases} \tag{22}$$

- If the priors of two classes are equal, i.e., $\pi_1 = \pi_2$, the Eq. (20) becomes:

$$2\left(\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)\right)^{\top}\boldsymbol{x} + \boldsymbol{\mu}_1^{\top}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^{\top}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_2 = 0, \tag{23}$$

whose left-hand-side expression can be considered as $\delta(\boldsymbol{x})$ in Eq. (22).

**Quadratic Discriminant Analysis**

# Quadratic Discriminant Analysis for Binary Classification

- In Quadratic Discriminant Analysis (QDA), we relax the assumption of equality of the covariance matrices:

$$\mathbf{\Sigma}_1 \neq \mathbf{\Sigma}_2, \qquad \ln(a^b) = b \ln a \qquad (24)$$

which means the covariances are not *necessarily* equal (if they are actually equal, the decision boundary will be linear and QDA reduces to LDA).

- Therefore, the Eq. (17) becomes:

$$\frac{1}{\sqrt{(2\pi)^d |\mathbf{\Sigma}_1|}} \exp\left(-\frac{(\mathbf{x} - \boldsymbol{\mu}_1)^\top \mathbf{\Sigma}_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)}{2}\right) \pi_1$$

$$= \frac{1}{\sqrt{(2\pi)^d |\mathbf{\Sigma}_2|}} \exp\left(-\frac{(\mathbf{x} - \boldsymbol{\mu}_2)^\top \mathbf{\Sigma}_2^{-1}(\mathbf{x} - \boldsymbol{\mu}_2)}{2}\right) \pi_2,$$

$$\overset{(a)}{\Longrightarrow} -\frac{d}{2}\ln(2\pi) - \frac{1}{2}\ln(|\mathbf{\Sigma}_1|) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^\top \mathbf{\Sigma}_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \ln(\pi_1)$$

$$= -\frac{d}{2}\ln(2\pi) - \frac{1}{2}\ln(|\mathbf{\Sigma}_2|) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^\top \mathbf{\Sigma}_2^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) + \ln(\pi_2),$$

where (a) takes natural logarithm from the sides of equation.

# Quadratic Discriminant Analysis for Binary Classification

- Recall Eq. (19):

$$(x - \mu_1)^\top \Sigma^{-1} (x - \mu_1) = x^\top \Sigma^{-1} x + \mu_1^\top \Sigma^{-1} \mu_1 - 2\mu_1^\top \Sigma^{-1} x.$$

- According to Eq. (19), we have:

$$-\left(\frac{1}{2}\right)\ln(|\Sigma_1|) - \left(\frac{1}{2}\right)x^\top \Sigma_1^{-1} x - \left(\frac{1}{2}\right)\mu_1^\top \Sigma_1^{-1} \mu_1 + \mu_1^\top \Sigma_1^{-1} x + \ln(\pi_1)$$
$$= -\left(\frac{1}{2}\right)\ln(|\Sigma_2|) - \left(\frac{1}{2}\right)x^\top \Sigma_2^{-1} x - \left(\frac{1}{2}\right)\mu_2^\top \Sigma_2^{-1} \mu_2 + \mu_2^\top \Sigma_2^{-1} x + \ln(\pi_2).$$

- Therefore, if we multiply the sides of equation by 2, we have:

$$x^\top (\Sigma_1 - \Sigma_2)^{-1} x + 2(\Sigma_2^{-1}\mu_2 - \Sigma_1^{-1}\mu_1)^\top x$$
$$+ (\mu_1^\top \Sigma_1^{-1} \mu_1 - \mu_2^\top \Sigma_2^{-1} \mu_2) + \ln\left(\frac{|\Sigma_1|}{|\Sigma_2|}\right) + 2\ln(\frac{\pi_2}{\pi_1}) = 0, \quad (25)$$

which is in the quadratic form $x^\top A x + b^\top x + c = 0.$

- Therefore, if we consider Gaussian distributions for the two classes, the decision boundary of classification is quadratic. Because of quadratic decision boundary which discriminates the two classes, this method is named *quadratic discriminant* analysis.

# Quadratic Discriminant Analysis for Binary Classification

- For obtaining Eq. (25), we brought the expressions to the right side which was corresponding to the second class; therefore, if we use $\delta(\boldsymbol{x}) : \mathbb{R}^d \to \mathbb{R}$ as the left-hand-side expression (function) in Eq. (25):

$$
\begin{aligned}
\delta(\boldsymbol{x}) := &\ \boldsymbol{x}^\top (\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2)^{-1} \boldsymbol{x} + 2\,(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1)^\top \boldsymbol{x} \\
&+ (\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) + \ln\!\Big(\frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_2|}\Big) + 2\,\ln\!\Big(\frac{\pi_2}{\pi_1}\Big),
\end{aligned}
\tag{26}
$$

the class of an instance $\boldsymbol{x}$ is estimated as the Eq. (22):

$$
\widehat{\mathcal{C}}(x) = \left\{ \begin{array}{ll} 1, & \text{if } \delta(\boldsymbol{x}) < 0, \\ 2, & \text{if } \delta(\boldsymbol{x}) > 0. \end{array} \right.
$$

- If the priors of two classes are equal, i.e., $\pi_1 = \pi_2$, the Eq. (20) becomes:

$$
\begin{aligned}
&\boldsymbol{x}^\top (\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2)^{-1} \boldsymbol{x} + 2\,(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1)^\top \boldsymbol{x} \\
&+ (\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) + \ln\!\Big(\frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_2|}\Big) = 0,
\end{aligned}
\tag{27}
$$

whose left-hand-side expression can be considered as $\delta(\boldsymbol{x})$ in Eq. (22).

if $\Sigma_1 = \Sigma_2$ ?

**LDA and QDA for Multi-class Classification**

# LDA and QDA for Multi-class Classification

- Now we consider underlined multiple classes, which can be more than two, indexed by $k \in \{1, \ldots, |\mathcal{C}|\}$. Recall Eq. (12) or (15) where we are using the scaled posterior, i.e., $f_k(\boldsymbol{x}) \pi_k$.
- According to Eq. (16), we have:

$$f_k(\boldsymbol{x}) \pi_k = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp\left(-\frac{(\boldsymbol{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_k)}{2}\right) \pi_k$$

Taking natural logarithm gives:

$$\ln(f_k(\boldsymbol{x}) \pi_k) = -\frac{d}{2}\ln(2\pi) - \frac{1}{2}\ln(|\boldsymbol{\Sigma}_k|) - \frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_k) + \ln(\pi_k).$$

- We drop the constant term $-(d/2)\ln(2\pi)$ which is the same for all classes (note that this term is multiplied before taking the logarithm). Thus, the scaled posterior of the $k$-th class becomes:

$$\delta_k(\boldsymbol{x}) := -\frac{1}{2}\ln(|\boldsymbol{\Sigma}_k|) - \frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_k) + \ln(\pi_k). \tag{28}$$

- In QDA, the class of the instance $\boldsymbol{x}$ is estimated as:

$$\widehat{\mathcal{C}}(\boldsymbol{x}) = \arg\max_k \delta_k(\boldsymbol{x}), \tag{29}$$

because it maximizes the posterior of that class. In this expression, $\delta_k(\boldsymbol{x})$ is Eq. (28).

# LDA and QDA for Multi-class Classification

- In LDA, we assume that the covariance matrices of the $k$ classes are equal:

$$\mathbf{\Sigma}_1 = \cdots = \mathbf{\Sigma}_{|\mathcal{C}|} = \mathbf{\Sigma}. \tag{30}$$

Therefore, the Eq. (28),

$$\delta_k(\mathbf{x}) := -\frac{1}{2}\ln(|\mathbf{\Sigma}_k|) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^\top \mathbf{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) + \ln(\pi_k),$$

becomes:

$$\delta_k(\mathbf{x}) = -\frac{1}{2}\ln(|\mathbf{\Sigma}|) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^\top \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) + \ln(\pi_k)$$

$$= -\frac{1}{2}\ln(|\mathbf{\Sigma}|) - \frac{1}{2}\mathbf{x}^\top \mathbf{\Sigma}^{-1}\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_k^\top \mathbf{\Sigma}^{-1}\boldsymbol{\mu}_k + \boldsymbol{\mu}_k^\top \mathbf{\Sigma}^{-1}\mathbf{x} + \ln(\pi_k).$$

- We drop the constant terms $-(1/2)\ln(|\mathbf{\Sigma}|)$ and $-(1/2)\,\mathbf{x}^\top \mathbf{\Sigma}^{-1}\mathbf{x}$ which are the same for all classes (note that before taking the logarithm, the term $-(1/2)\ln(|\mathbf{\Sigma}|)$ is multiplied and the term $-(1/2)\,\mathbf{x}^\top \mathbf{\Sigma}^{-1}\mathbf{x}$ is multiplied as an exponential term). Thus, the scaled posterior of the $k$-th class becomes:

$$\delta_k(\mathbf{x}) := \boldsymbol{\mu}_k^\top \mathbf{\Sigma}^{-1}\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_k^\top \mathbf{\Sigma}^{-1}\boldsymbol{\mu}_k + \ln(\pi_k). \tag{31}$$

- In LDA, the class of the instance $\mathbf{x}$ is determined by Eq. (29), where $\delta(\mathbf{x})$ is Eq. (31), because it maximizes the posterior of that class.

**Estimation of Parameters in LDA and QDA**

# Estimation of Parameters in LDA and QDA

- In LDA and QDA, we have several parameters which are required in order to calculate the posteriors. These parameters are the **means** and the **covariance matrices** of classes and the **priors** of classes.

- The priors of the classes are very tricky to calculate. It is somewhat a chicken and egg problem because we want to know the class probabilities (priors) to estimate the class of an instance but we do not have the priors and should estimate them.

- Usually, the **prior** of the $k$-th class is estimated according to the sample size of the $k$-th class:

$$\widehat{\pi}_k = \frac{n_k}{n}, \tag{32}$$

where $n_k$ and $n$ are the number of training instances in the $k$-th class and in total, respectively. This estimation considers Bernoulli distribution for choosing every instance out of the overall training set to be in the $k$-th class.

# Estimation of Parameters in LDA and QDA

- The **mean** of the $k$-th class can be estimated using the Maximum Likelihood Estimation (MLE), or Method of Moments (MOM), for the mean of a Gaussian distribution:

$$\mathbb{R}^d \ni \widehat{\boldsymbol{\mu}}_k = \frac{1}{n_k} \sum_{i=1}^{n} \mathbf{x}_i \mathbb{I}(\mathcal{C}(\mathbf{x}_i) = k), \tag{33}$$

where $\mathbb{I}(.)$ is the indicator function which is one and zero if its condition is satisfied and not satisfied, respectively.

- In QDA, the covariance matrix of the $k$-th class is estimated using MLE:

$$\mathbb{R}^{d \times d} \ni \widehat{\boldsymbol{\Sigma}}_k = \frac{1}{n_k} \sum_{i=1}^{n} (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_k)^\top \mathbb{I}(\mathcal{C}(\mathbf{x}_i) = k). \tag{34}$$

Or we can use the *unbiased* estimation of the covariance matrix:

$$E(\hat{\Sigma}) = \frac{n_k - 1}{n_k} \Sigma$$

$$\mathbb{R}^{d \times d} \ni \widehat{\boldsymbol{\Sigma}}_k = \frac{1}{n_k - 1} \sum_{i=1}^{n} (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_k)^\top \mathbb{I}(\mathcal{C}(\mathbf{x}_i) = k). \tag{35}$$

- In LDA, we assume that the covariance matrices of the classes are equal; therefore, we use the weighted average of the estimated covariance matrices as the common covariance matrix in LDA:

$$\mathbb{R}^{d \times d} \ni \widehat{\boldsymbol{\Sigma}} = \frac{\sum_{k=1}^{|\mathcal{C}|} n_k \widehat{\boldsymbol{\Sigma}}_k}{\sum_{r=1}^{|\mathcal{C}|} n_r} = \frac{\sum_{k=1}^{|\mathcal{C}|} n_k \widehat{\boldsymbol{\Sigma}}_k}{n}, \tag{36}$$

where the weights are the cardinality of the classes.

$$P(\cancel{} c \mid X) = \frac{P(X \mid c) \, P(c)}{\cancel{P(X)}}$$

posterior

likelihood
(class conditional)

prior
of class

prior
of data

marginal

$$P(X \mid c) \, P(c) = P(X, c)$$

$$P(X) = \sum_c P(x, c)$$
$$= \sum_c P(X \mid c) \, P(c)$$

# Examples for LDA and QDA
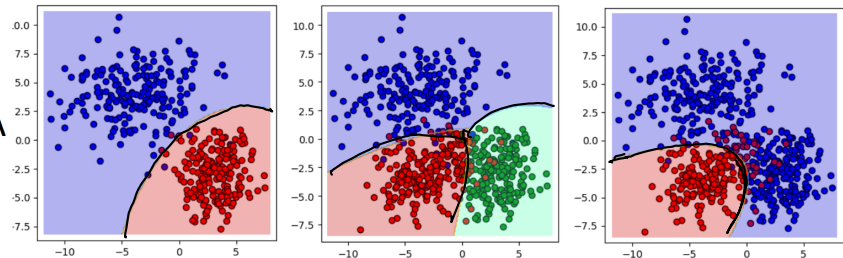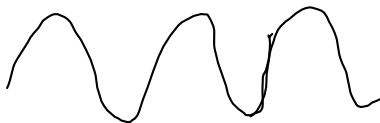
Code of these plots and example in my GitHub:
`https://github.com/bghojogh/Linear-Quadratic-Discriminant-Analysis`

Gabor filter



LDA and QDA are
Metric Learning!

central limit theorem



Poisson distribution

# LDA and QDA are Metric Learning!

- Recall Eq. (28) which is the scaled posterior for the QDA:

$$\delta_k(\boldsymbol{x}) := -\frac{1}{2}\ln(|\boldsymbol{\Sigma}_k|) - \frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_k) + \ln(\pi_k).$$

- First, assume that the covariance matrices are all equal (as we have in LDA) and they all are the identity matrix:

$$\boldsymbol{\Sigma}_1 = \cdots = \boldsymbol{\Sigma}_{|\mathcal{C}|} = \boldsymbol{I}, \tag{37}$$

which means that all the classes are assumed to be spherically distributed in the $d$ dimensional space. After this assumption, the Eq. (28) becomes:

$$\delta_k(\boldsymbol{x}) = -\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_k)^\top(\boldsymbol{x} - \boldsymbol{\mu}_k) + \ln(\pi_k), \tag{38}$$

because $|\boldsymbol{I}| = 1$, $\ln(1) = 0$, and $\boldsymbol{I}^{-1} = \boldsymbol{I}$.

- If we assume that the priors are all equal, the term $\ln(\pi_k)$ is constant and can be dropped:

$$\delta_k(\boldsymbol{x}) = -\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_k)^\top(\boldsymbol{x} - \boldsymbol{\mu}_k) = -\frac{1}{2}d_k^2, \tag{39}$$

where $d_k$ is the Euclidean distance from the mean of the $k$-th class:

$$d_k = ||\boldsymbol{x} - \boldsymbol{\mu}_k||_2 = \sqrt{(\boldsymbol{x} - \boldsymbol{\mu}_k)^\top(\boldsymbol{x} - \boldsymbol{\mu}_k)}. \tag{40}$$

- Thus, the QDA or LDA reduce to simple Euclidean distance from the means of classes if the covariance matrices are all identity matrix and the priors are equal. Simple distance from the mean of classes is one of the simplest classification methods where the used metric is Euclidean distance.
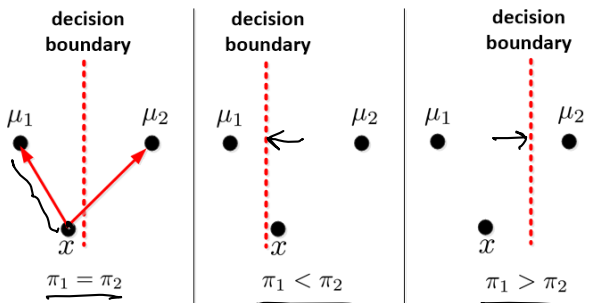
# LDA and QDA are Metric Learning!

- Now, consider the case where still the **covariance matrices are all identity matrix** but the **priors are not equal**. In this case, we have Eq. (38):

$$\delta_k(x) = -\frac{1}{2}(x - \mu_k)^\top (x - \mu_k) + \ln(\pi_k).$$

- If we take an exponential (inverse of logarithm) from this expression, the $\pi_k$ becomes a scale factor (weight). This means that we still are using distance metric to measure the distance of an instance from the means of classes but we are **scaling the distances by the priors of classes**. If a class happens more, i.e., its **prior is larger**, it must have a larger posterior so we **reduce the distance** from the mean of its class. In other words, we **move the decision boundary** according to the prior of classes.

# LDA and QDA are Metric Learning!

- As the next step, consider a more general case where the **covariance matrices are not equal** as we have in **QDA**.
- We apply Singular Value Decomposition (SVD) to the covariance matrix of the $k$-th class:

$$\boldsymbol{\Sigma}_k = \boldsymbol{U}_k \, \boldsymbol{\Lambda}_k \, \boldsymbol{U}_k^\top, \quad \longleftarrow$$

where the left and right matrices of singular vectors are equal because the covariance matrix is symmetric. Therefore:

$$\boldsymbol{\Sigma}_k^{-1} = \boldsymbol{U}_k \, \boldsymbol{\Lambda}_k^{-1} \, \boldsymbol{U}_k^\top,$$

$$U^{-T} = U^{T^{-1}} = U$$

where $\boldsymbol{U}_k^{-1} = \boldsymbol{U}_k^\top$ because it is an orthogonal matrix.

- Therefore, we can simplify the following term:

$$(\boldsymbol{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\boldsymbol{x} - \boldsymbol{\mu}_k) = (\boldsymbol{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{U}_k \boldsymbol{\Lambda}_k^{-1} \boldsymbol{U}_k^\top (\boldsymbol{x} - \boldsymbol{\mu}_k)$$

$$= (\boldsymbol{U}_k^\top \boldsymbol{x} - \boldsymbol{U}_k^\top \boldsymbol{\mu}_k) \boldsymbol{\Lambda}_k^{-1} (\boldsymbol{U}_k^\top \boldsymbol{x} - \boldsymbol{U}_k^\top \boldsymbol{\mu}_k).$$

$$(x^\top U - \mu^\top U) = (x^\top - \mu^\top) U = (x - \mu)^\top U$$

# LDA and QDA are Metric Learning!

$A = A^{1/2} A^{1/2}$

- We found:

$$(x - \mu_k)^\top \Sigma_k^{-1} (x - \mu_k) = (U_k^\top x - U_k^\top \mu_k)^\top \Lambda_k^{-1} U_k^\top x - U_k^\top \mu_k).$$

- As $\Lambda_k^{-1}$ is a diagonal matrix with non-negative elements (because it is covariance), we can decompose it as:

$$\Lambda_k^{-1} = \Lambda_k^{-1/2} \Lambda_k^{-1/2}.$$

Therefore:

$$(U_k^\top x - U_k^\top \mu_k)^\top \Lambda_k^{-1} (U_k^\top x - U_k^\top \mu_k) = (U_k^\top x - U_k^\top \mu_k)^\top \Lambda_k^{-1/2} \Lambda_k^{-1/2} (U_k^\top x - U_k^\top \mu_k)$$

$$\overset{(a)}{=} (\Lambda_k^{-1/2} U_k^\top x - \Lambda_k^{-1/2} U_k^\top \mu_k)^\top (\Lambda_k^{-1/2} U_k^\top x - \Lambda_k^{-1/2} U_k^\top \mu_k),$$

$= \left[ \phi_k(x) - \phi_k(\mu_k) \right]^\top \left( \phi_k(x) - \phi_k(\mu_k) \right)$

where $(a)$ is because $\Lambda_k^{-\top/2} = \Lambda_k^{-1/2}$ because it is diagonal.

- We define the following transformation:

$$\phi_k : x \mapsto \Lambda_k^{-1/2} U_k^\top x, \qquad (41)$$

which also results in the transformation of the mean: $\phi_k : \mu \mapsto \Lambda_k^{-1/2} U_k^\top \mu$.

# LDA and QDA are Metric Learning!

- We had:

$$\underbrace{\phi_k} : \boldsymbol{x} \mapsto \boldsymbol{\Lambda}_k^{-1/2} \boldsymbol{U}_k^\top \boldsymbol{x}.$$

- Therefore, the Eq. (28),

$$\delta_k(\boldsymbol{x}) := -\frac{1}{2} \ln(|\boldsymbol{\Sigma}_k|) - \frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_k) + \ln(\pi_k),$$

  can be restated as:

$$\delta_k(\boldsymbol{x}) = -\frac{1}{2} \ln(|\boldsymbol{\Sigma}_k|) - \frac{1}{2}\big(\phi_k(\boldsymbol{x}) - \phi_k(\boldsymbol{\mu}_k)\big)^\top \big(\phi_k(\boldsymbol{x}) - \phi_k(\boldsymbol{\mu}_k)\big) + \ln(\pi_k). \tag{42}$$

- Ignoring the terms $-(1/2)\ln(|\boldsymbol{\Sigma}_k|)$ and $\ln(\pi_k)$, we can see that the **transformation** has **changed the covariance matrix of the class to identity matrix**. Therefore, the QDA (and also LDA) can be seen as **simple comparison of distances from the means of classes after applying a transformation to the data of every class**. In other words, we are **learning the metric using the SVD of covariance matrix of every class**. Thus, LDA and QDA can be seen as **metric learning** [2, 3] in a perspective.

# LDA and QDA are Metric Learning!



$$(x - \mu)^\top (x - \mu)$$

$$(x - \mu)^\top \Sigma^{-1} (x - \mu)$$

- Note that in metric learning, a valid distance metric is defined as [2]:

$$d_A^2(x, \mu_k) := \|x - \mu_k\|_A^2 = (x - \mu_k)^\top A (x - \mu_k), \qquad (43)$$
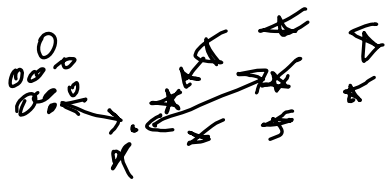
  where $A$ is a positive semi-definite matrix, i.e., $A \succeq 0$.
- In QDA, we are also using $(x - \mu_k)^\top \Sigma_k^{-1} (x - \mu_k)$.
- The covariance matrix is positive semi-definite according to the characteristics of covariance matrix. Moreover, according to characteristics of a positive semi-definite matrix, the inverse of a positive semi-definite matrix is positive semi-definite so $\Sigma_k^{-1} \succeq 0$.
- Therefore, QDA is using metric learning.

$$(x - \mu)^\top A (x - \mu) = (x - \mu)^\top U U^\top (x - \mu)$$

$$A = V \Lambda V^\top = \underbrace{V \Lambda^{\frac{1}{2}}}_{U} \Lambda^{\frac{1}{2}} V^\top = U U^\top$$

$$(U^\top x - U^\top \mu)^\top (U^\top x - U^\top \mu)$$

# Acknowledgment

# References

[1] T. Hastie, R. Tibshirani, J. H. Friedman, and J. H. Friedman, *The elements of statistical learning: data mining, inference, and prediction*, vol. 2. Springer, 2009.

[2] L. Yang and R. Jin, "Distance metric learning: A comprehensive survey," tech. rep., Department of Computer Science and Engineering, Michigan State University, 2006.

[3] B. Kulis, "Metric learning: A survey," *Foundations and Trends® in Machine Learning*, vol. 5, no. 4, pp. 287–364, 2013.

[4] B. Ghojogh and M. Crowley, "Linear and quadratic discriminant analysis: Tutorial," *arXiv preprint arXiv:1906.02590*, 2019.