



1. **EXAMple:** The lecture introduced Fisher's linear discriminant analysis. Fisher's LDA is a relatively simple algorithm, but there is an even simpler one: To find a projection $\mathbf{w} \in \mathbb{R}^d$ that "separates" a dataset $X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ into two classes $\mathcal{C}_1, \mathcal{C}_2$ (i.e. index sets with $\mathcal{C}_1 \cup \mathcal{C}_2 = [1, \dots, N]$ and $\mathcal{C}_1 \cap \mathcal{C}_2 = \emptyset$), we can consider the class means

$$\mathbf{m}_i := \frac{1}{N_i} \sum_{n \in \mathcal{C}_i} \mathbf{x}_n \quad (\text{where } N_i := |\mathcal{C}_i|),$$

and choose a projection $\mathbf{w} \in \mathbb{R}^D$ such that the distance of projected classes $m_1 - m_2 = \mathbf{w}^\top(\mathbf{m}_1 - \mathbf{m}_2)$ is maximized.¹ Of course, this can be made arbitrarily large by increasing $\|\mathbf{w}\|$, so we have to constrain $\|\mathbf{w}\| = 1$ with a Lagrange multiplier λ : Find this "optimal" projection \mathbf{w} by computing the gradient of

$$u(\mathbf{w}) = m_1 - m_2 + \lambda(1 - \|\mathbf{w}\|^2),$$

setting it to zero, and using the constraint $\|\mathbf{w}\| = 1$ to find the value of λ .

2. **Theory Question:** Fisher's choice of univariate projection for linear discriminant analysis was introduced in the lecture. However, the actual *discriminant* itself was not derived. This latter step is the goal of this exercise. Consider data as in Exercise 1 above, with class means \mathbf{m}_i as above and unnormalized covariances

$$S_i := \sum_{n \in \mathcal{C}_i} (\mathbf{x}_n - \mathbf{m}_i)(\mathbf{x}_n - \mathbf{m}_i)^\top.$$

LDA finds a projection $\mathbf{w} \in \mathbb{R}^D$. The projected points $\mathbf{x}_n = \mathbf{w}^\top \mathbf{x}_n$ then have means $m_i = \mathbf{w}^\top \mathbf{m}_i$ and variance $s_i^2 = \frac{1}{N_i} \mathbf{w}^\top S_i \mathbf{w}$. We can assume² that this defines two separate Gaussian distributions $p(x_n | n \in \mathcal{C}_i) = \mathcal{N}(x_n; m_i, s_i^2)$. We can then consider a generative Gaussian mixture model that assigns prior probability $p(n \in \mathcal{C}_i) = N_i/(N_1 + N_2)$ to each class, then draws x_n by first drawing class i from $p(n \in \mathcal{C}_i)$, then drawing x_n from $p(x_n | n \in \mathcal{C}_i)$.

One way to motivate a discriminant (i.e. a rule to predict class 1 or 2 from x_n) is thus to predict, at x_n , the class with higher posterior probability. The discriminant then lies at the point where the log odds cancel, i.e. where

$$\log \frac{p(\mathcal{C}_1 | x_n)}{p(\mathcal{C}_2 | x_n)} = \log \frac{p(x_n | n \in \mathcal{C}_1)p(n \in \mathcal{C}_1)}{p(x_n | n \in \mathcal{C}_2)p(n \in \mathcal{C}_2)} = 0 \quad (1)$$

- (a) Solve Equation 1 above for $x \in \mathbb{R}$, by plugging in the definitions of the terms from above.
 - (b) You will find that you arrive at a quadratic equation with two solutions $x_{1,2}$. Which one of these two is the discriminant?
3. **Practical Question:** You can find this week's practical exercise in `Exercise_09.ipynb`

¹This reason that this discriminant is not used in practice is that it does not necessarily separate the classes well, unless the the data covariance aligns with the vector $\mathbf{m}_1 - \mathbf{m}_2$.

²A frequently cited motivation for this assumption is that each scalar $x_n = \mathbf{w}^\top \mathbf{x}_n$ is a sum of d random variables. Thus, the Central Limit Theorem suggests that the x_n are indeed approximately Gaussian.