

Ex 4.1

W = Within class Variance - How much each observation varies from its class mean.

$$W = \sum_K \sum_{g \in K} (x_i - \hat{\mu}_K)^T (x_i - \hat{\mu}_K) / (N - K)$$

B = Between class Variance - How much does each cluster vary from the overall mean, $\hat{\mu} = \sum_K \hat{\pi}_K \hat{\mu}_K$

$$B = \sum_K \hat{\pi}_K (\hat{\mu}_K - \hat{\mu})^T (\hat{\mu}_K - \hat{\mu})$$

Now K centroids in P dimensional input space lie in an affine subspace of dimension $\leq K-1$. If P is much larger than K , this is a considerable drop in dimension.

Also notice, to find the nearest centroid to a point we need only to project it to this affine subspace and find the nearest centroid from that point.

This is the inherent dimensionality reduction in LDA. To determine which dimension subspace ($\leq K-1$) was optimal for LDA, Fisher defined that to mean the linear combination $Z = a^T X$ such that between class variance is maximised relative to within class variance.

\Rightarrow If $Z = a^T X$ is some linear combination of X then,

$$\text{Var}_B(Z) = \text{Var}_B(a^T X) = a^T \text{Var}_B(X) a = a^T B a$$

$$\text{and } \text{Var}_W(Z) = \dots = a^T W a$$

Thus the Problem amounts to $\max_a \frac{a^T B a}{a^T W a}$ (4.15)

or equivalently $\max_a a^T B a$ subject to $a^T W a = 1$

e.g. a constrained maximisation Problem solvable by Lagrange Multipliers.

$\Rightarrow \max_a a^T B a$ subject to $a^T W a - 1 = 0$

$$L(a) = a^T B a - \lambda(a^T W a - 1)$$

$$\nabla_a L(a) = \nabla_a(a^T B a) - \lambda \nabla_a(a^T W a) = 0$$

$$\Rightarrow \nabla_a(a^T B a) = \lambda \nabla_a(a^T W a)$$

$$= 2 B a = 2 \lambda W a$$

$$\Rightarrow B a = \lambda W a$$

$\Rightarrow W^{-1} B a = \lambda a$ which is just a standard eigenvalue Problem.

The eigenvector (a) corresponding to the largest eigenvalue (λ) will maximise the expression:

$$\frac{a^T B a}{a^T W a} = \frac{a^T (\lambda W a)}{a^T W a} = \lambda$$