

Dimensionality Reduction through LDA

- We studied **dimensionality reduction** through PCA.
- LDA can be used for dimensionality reduction.
It can be considered as a “supervised PCA”.
- Given $\mathcal{D} = \mathcal{D}_1 \sqcup \mathcal{D}_2 \sqcup \dots \sqcup \mathcal{D}_s$ (disjoint union),
set

$$N_t := \#(\mathcal{D}_t), \quad t = 1, 2, \dots, s, \quad N := N_1 + \dots + N_s.$$

- Main Idea: Shrink class \mathcal{D}_t into a single point $\boldsymbol{\mu}_t$ and do PCA.
 \rightsquigarrow Principal directions for **classes**.

Choose

$$\boldsymbol{\mu}_t = \frac{1}{N_t} \sum_{\mathbf{x} \in \mathcal{D}_t} \mathbf{x} \quad \text{and} \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{\mathbf{x} \in \mathcal{D}} \mathbf{x}.$$

However, the point $\boldsymbol{\mu}_t$ must have multiplicity N_t .

- Find a feature vector \mathbf{a} which maximizes

$$\mathbf{a}^\top \left(\sum_{t=1}^s N_t (\boldsymbol{\mu}_t - \boldsymbol{\mu})(\boldsymbol{\mu}_t - \boldsymbol{\mu})^\top \right) \mathbf{a}.$$

- A feature vector \mathbf{a} may capture more variability from one class than from the other classes.

To normalize, we fix the sum of variances. That is,

$$\sum_{t=1}^s \mathbf{a}^\top \left(\sum_{\mathbf{x} \in \mathcal{D}_t} (\mathbf{x} - \boldsymbol{\mu}_t)(\mathbf{x} - \boldsymbol{\mu}_t)^\top \right) \mathbf{a} = 1.$$

- Define

$$B = \sum_{t=1}^s N_t (\boldsymbol{\mu}_t - \boldsymbol{\mu})(\boldsymbol{\mu}_t - \boldsymbol{\mu})^\top$$
$$W = \sum_{t=1}^s \sum_{\mathbf{x} \in \mathcal{D}_t} (\mathbf{x} - \boldsymbol{\mu}_t)(\mathbf{x} - \boldsymbol{\mu}_t)^\top.$$

- Task: Maximize $\mathbf{a}^\top B \mathbf{a}$, subject to $\mathbf{a}^\top W \mathbf{a} = 1$.

- Lagrange Multiplier

$$F = \mathbf{a}^\top B \mathbf{a} - \lambda (\mathbf{a}^\top W \mathbf{a} - 1)$$

- We obtain

$$\nabla F = 2B\mathbf{a} - 2\lambda W\mathbf{a} = 0 \iff W^{-1}B\mathbf{a} = \lambda\mathbf{a}$$

Critical points are eigenvectors of $W^{-1}B$.

Note that $W^{-1}B$ is symmetric.

- If \mathbf{a} is an eigenvector of $W^{-1}B$ such that $\mathbf{a}^\top W\mathbf{a} = 1$, then

$$\mathbf{a}^\top B\mathbf{a} = \mathbf{a}^\top WW^{-1}B\mathbf{a} = \lambda\mathbf{a}^\top W\mathbf{a} = \lambda.$$

Thus an eigenvector of the **largest eigenvalue** is the first principal direction.

- Take k' -many principal directions for $k' < k$.

Project data points onto the subspace of the principal directions.

~~~~~> Dimensionality Reduction